

Supplementary materials for “A Study on Truncated Newton Methods for Linear Classification”

Leonardo Galli

*Department of Information Engineering
University of Florence
Via Santa Marta 3, Firenze, Italia*

LEONARDO.GALLI@UNIFI.IT

Chih-Jen Lin

*Department of Computer Science
National Taiwan University
Taipei, Taiwan, 106*

CJLIN@CSIE.NTU.EDU.TW

We report here some equations from the original paper that might be useful in the rest of the supplementary materials file.

The **linear classification problem**:

$$\min_{\mathbf{w}} f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \xi(y_i \mathbf{w}^T \mathbf{x}_i), \quad (\text{i})$$

where (y_i, \mathbf{x}_i) , $i = 1, \dots, l$ are the training data, $y_i = \pm 1$ is the label, $\mathbf{x}_i \in \mathbb{R}^n$ is a feature vector, $\mathbf{w}^T \mathbf{w} / 2$ is the regularization term, $C > 0$ is a regularization parameter and $\xi(y_i \mathbf{w}^T \mathbf{x}_i)$ is a generic LC^1 (continuously differentiable with locally Lipschitz continuous gradient) convex loss function. With $\mathbf{w}^T \mathbf{w}$, f is a LC^1 strongly convex function and the minimum \mathbf{w}^* of f exists and is unique.

The following two **losses** respectively correspond to logistic regression (C^2 , i.e. twice continuously differentiable) and l2-loss linear SVM (LC^1)

$$\begin{aligned} \xi_{\text{LR}}(y \mathbf{w}^T \mathbf{x}) &= \log(1 + \exp(-y \mathbf{w}^T \mathbf{x})) \\ \xi_{\text{L2}}(y \mathbf{w}^T \mathbf{x}) &= (\max(0, 1 - y \mathbf{w}^T \mathbf{x}))^2. \end{aligned} \quad (\text{ii})$$

The **gradient** and the **Hessian** of the objective function f :

$$\begin{aligned} \mathbf{g} &:= \nabla f(\mathbf{w}) = \mathbf{w} + C \sum_{i=1}^l \xi'(y_i \mathbf{w}^T \mathbf{x}_i) y_i \mathbf{x}_i, \\ H &:= \nabla^2 f(\mathbf{w}) = I + C X^T D X, \end{aligned} \quad (\text{iii})$$

where I is the identity matrix, $X = [\mathbf{x}_1, \dots, \mathbf{x}_l]^T$ is the data matrix and D is a diagonal matrix with $D_{ii} = \xi''(y_i \mathbf{w}^T \mathbf{x}_i)$. The **Newton equation**:

$$\nabla^2 f(\mathbf{w}_k) \mathbf{s} = -\nabla f(\mathbf{w}_k). \quad (\text{iv})$$

Generic **truncation criterion**:

$$\text{ratio}(\mathbf{s}_k^j) \leq \eta_k, \quad (\text{v})$$

Armijo condition:

$$f(\mathbf{w}_k + \omega_k \mathbf{s}_k) \leq f(\mathbf{w}_k) + \gamma \omega_k \mathbf{g}_k^T \mathbf{s}_k, \quad \text{with } 0 < \gamma < 1. \quad (\text{vi})$$

Algorithm I: Conjugate Gradient

Input: $\mathbf{s}_k^1 = \mathbf{0} \in \mathbb{R}^n$, $\mathbf{d}_k^1 = -\nabla Q_1 = -\mathbf{g}_k$
1 for $j=1, 2, \dots$ **do**
2 $\alpha_k^j = \frac{\|\nabla Q_j\|^2}{\mathbf{d}_k^{jT} H_k \mathbf{d}_k^j}$
3 $\mathbf{s}_k^{j+1} = \mathbf{s}_k^j + \alpha_k^j \mathbf{d}_k^j$
4 $\nabla Q_{j+1} = \nabla Q_j + \alpha_k^j H_k \mathbf{d}_k^j$
5 **if** (v) *is satisfied* **then**
6 $\mathbf{s}_k = \mathbf{s}_k^{j+1}$
7 **return** \mathbf{s}_k
8 $\beta_j = \frac{\|\nabla Q_{j+1}\|^2}{\|\nabla Q_j\|^2}$
9 $\mathbf{d}_k^{j+1} = -\nabla Q_{j+1} + \beta_j \mathbf{d}_k^j$

Wolfe condition:

$$\nabla f(\mathbf{w}_k + \omega_k \mathbf{s}_k)^T \mathbf{s}_k \geq \gamma \mathbf{g}_k^T \mathbf{s}_k, \quad \text{with } 0 < \gamma < 1. \quad (\text{vii})$$

Level set:

$$\mathcal{L}_1 := \{\mathbf{w} \in \mathbb{R}^n : f(\mathbf{w}) \leq f(\mathbf{w}_1)\}. \quad (\text{viii})$$

This lemma can be obtained directly from the structure of the losses and it will be needed in Section I.

Lemma I *The Hessian matrix (iii) is bounded, i.e. there exist two constants $M_1, M_2 > 0$ such that*

$$M_1 \geq \|H_k\| \geq M_2 \quad \forall k. \quad (\text{ix})$$

Proof By the definition of the losses functions (ii) we have that $0 \leq D_{ii} \leq 2$ (see Appendix A for details). Thus from (iii) we have that $\forall k$ H_k is bounded. Furthermore, with the matrix I in (iii), H_k is bounded above zero. \blacksquare

I. Global and Local Convergence of TNCG

In this section we relax f to be any $f \in \text{LC}^1$.

I.1 Global Convergence by Treating TNCG as a Common-Directions Algorithm

In Algorithm I we detail the instructions of the conjugate gradient method for solving a generic quadratic strongly convex function of the shape $Q(\mathbf{s}) = \mathbf{g}_k^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T H_k \mathbf{s}$. We define $\nabla Q_j := \nabla Q(\mathbf{s}_k^j)$, where $\nabla Q(\mathbf{s}) = \mathbf{g}_k + H_k \mathbf{s}$.

Lemma II *Let \mathbf{s}_k be the direction returned by Algorithm I and assume that it terminates after m_k iterations. Then $\mathbf{s}_k = \sum_{j=1}^{m_k} \alpha_k^j \mathbf{d}_k^j$, where $P_k = \{\mathbf{d}_k^1, \dots, \mathbf{d}_k^{m_k}\}$ are conjugate directions*

Algorithm II: Truncated Newton Conjugate Gradient

Input: $\mathbf{w}_1 \in \mathbb{R}^n$ starting point

- 1 **for** $k = 1, 2, \dots$ **do**
- 2 compute the direction \mathbf{s}_k by approximately solving (iv) with a CG method, until (v) is satisfied
- 3 compute a step length ω_k
- 4 $\mathbf{w}_{k+1} = \mathbf{w}_k + \omega_k \mathbf{s}_k$

with respect to H_k and $\{\alpha_k^1, \dots, \alpha_k^{m_k}\}$ are scalars defined to minimize

$$\begin{aligned} \min_{\boldsymbol{\alpha}_k} \mathbf{g}_k^T \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T H_k \mathbf{s}_k \\ \text{s.t. } \mathbf{s}_k = P_k \boldsymbol{\alpha}_k, \end{aligned} \tag{x}$$

where $P_k := [\mathbf{d}_k^1 | \dots | \mathbf{d}_k^{m_k}] \in \mathbb{R}^{n \times m}$ and $\mathbf{g}_k \in \{\mathbf{d}_k^1, \dots, \mathbf{d}_k^{m_k}\}$.

Proof From instructions of Algorithm I we directly get that

$$\mathbf{s}_k = \mathbf{s}_k^{m_k+1} = \sum_{j=1}^{m_k} \alpha_k^j \mathbf{d}_k^j.$$

In addition, from Theorem 5.2 of Nocedal and Wright (1999) and the fact that $\mathbf{s}_k^1 = \mathbf{0}$ we have that \mathbf{s}_k is the minimizer of (x) on the span of the conjugate directions $\{\mathbf{d}_k^1, \dots, \mathbf{d}_k^{m_k}\}$. This concludes the proof. \blacksquare

We now point out that in the framework (Lee et al., 2017), the total number of common-directions m is assumed to be fixed through the whole optimization procedure. Since the termination criterion (v) is employed in Algorithm I, in the TNCG case this number m_k depends on the iteration k . On the other hand, it is possible to see that in the proofs of Lee et al. (2017), m is never really required to be fixed. In fact, as long as m is bounded, the problem (x) is well defined and all the properties required for global convergence follow from the fact that \mathbf{w}_k is its solution, together with the fact that $\mathbf{g}_k \in \{\mathbf{d}_k^1, \dots, \mathbf{d}_k^{m_k}\}$. Note that problem (x) is highly dependent on k , and the dimensionality of $\boldsymbol{\alpha}_k$ can be seen as a detail of the internal procedure for obtaining \mathbf{w}_k . For these reasons, we can relax the assumption that m is a fixed constant, and we instead require that it is bounded. In the case of Algorithm I, Theorem 5.1 from Nocedal and Wright (1999) ensures that the procedure terminates within n iterations, so we have that m_k is always bounded by n .

1.2 Global Convergence Following a More Classical Path

We first need to prove that the direction \mathbf{s}_k obtained by Algorithm I is able to satisfy (xi) and (xii). In this section we relax the assumption that H_k employed in $Q(\mathbf{s})$ of Algorithm I is the Hessian (or the generalized Hessian) of f . In fact, it is interesting to point out that for obtaining global convergence of $\{\mathbf{w}_k\}$, H_k can be a generic positive definite matrix. Indeed in general, second order information of f are not needed to obtain global convergence to a stationary point.

Proposition III Let $f \in LC^1$. In addition, assume that H_k is positive definite and $\forall k H_k$ is bounded as in (ix). Let \mathbf{s}_k be generated by Algorithm I. Then there exist two constants $a_1 > 0$ and $a_2 > 0$ such that

$$\nabla f(\mathbf{w}_k)^T \mathbf{s}_k \leq -a_1 \|\nabla f(\mathbf{w}_k)\|^2, \quad (\text{xi})$$

$$\|\mathbf{s}_k\| \leq a_2 \|\nabla f(\mathbf{w}_k)\|. \quad (\text{xii})$$

Proof If we pre-multiply

$$\nabla Q_j = \nabla Q_1 + \sum_{i=1}^{j-1} \alpha_i H_k \mathbf{d}_k^i$$

by \mathbf{d}_k^j we obtain

$$\mathbf{d}_k^{jT} \nabla Q_j = \mathbf{d}_k^{jT} \nabla Q_1, \quad (\text{xiii})$$

since \mathbf{d}_k^j are all mutually conjugated. From Theorem 5.1 of Hestenes and Stiefel (1952) we have that

$$\|\nabla Q_j\|^2 = -\nabla Q_i^T \mathbf{d}_k^j \quad \forall 1 \leq i \leq j, \quad (\text{xiv})$$

which together with (xiii) let us obtain that

$$\alpha^j = \frac{\|\nabla Q_j\|^2}{\mathbf{d}_k^{jT} H_k \mathbf{d}_k^j} = -\frac{\nabla Q_j^T \mathbf{d}_k^j}{\mathbf{d}_k^{jT} H_k \mathbf{d}_k^j} = -\frac{\nabla Q_1^T \mathbf{d}_k^j}{\mathbf{d}_k^{jT} H_k \mathbf{d}_k^j}. \quad (\text{xv})$$

From (xv) we get that

$$\mathbf{s}_k = \mathbf{s}_k^{j+1} = \sum_{i=1}^j \alpha^i \mathbf{d}_k^i = -\sum_{i=1}^j \frac{\nabla Q_1^T \mathbf{d}_k^i}{\mathbf{d}_k^{iT} H_k \mathbf{d}_k^i} \mathbf{d}_k^i. \quad (\text{xvi})$$

Thus,

$$\begin{aligned} \nabla f(\mathbf{w}_k)^T \mathbf{s}_k &= -\sum_{i=1}^j \frac{\left(\nabla f(\mathbf{w}_k)^T \mathbf{d}_k^i\right)^2}{\mathbf{d}_k^{iT} H_k \mathbf{d}_k^i} \\ &\leq -\frac{\left(\nabla f(\mathbf{w}_k)^T \nabla f(\mathbf{w}_k)\right)^2}{\mathbf{d}_k^{1T} H_k \mathbf{d}_k^1}, \end{aligned}$$

where the equality follows from the fact that $\nabla Q_1 = \nabla f(\mathbf{w}_k)$ and the inequality follows from the fact that H_k is positive definite and from the fact that in the algorithm initialization we have $\mathbf{d}_k^1 = \nabla f(\mathbf{w}_k)$. Thus, \mathbf{s}_k is a descent direction and we get

$$|\nabla f(\mathbf{w}_k)^T \mathbf{s}_k| \geq \frac{\|\nabla f(\mathbf{w}_k)\|^4}{\|\nabla f(\mathbf{w}_k)\|^2 \|H_k\|} \geq \frac{1}{M_1} \|\nabla f(\mathbf{w}_k)\|^2,$$

where the last inequality follows from upper-boundedness of H with the constant $M_1 > 0$. This proves that (xi) is satisfied with $a_1 = \frac{1}{M_1}$. Now to prove (xii) we first observe that H_k is also lower bounded, i.e.

$$\mathbf{d}_k^{jT} H_k \mathbf{d}_k^j \geq M_2 \mathbf{d}_k^{jT} I \mathbf{d}_k^j \geq M_2 \|\mathbf{d}_k^j\|^2, \quad (\text{xvii})$$

where I is the identity matrix. Thus, from (xvi) and (xvii) we can obtain

$$\begin{aligned} \|\mathbf{s}_k\| &= \|\mathbf{s}_k^{j+1}\| = \sum_{i=1}^j \left| \frac{\nabla Q_1^T \mathbf{d}_k^i}{\mathbf{d}_k^{iT} H_k \mathbf{d}_k^i} \right| \|\mathbf{d}_k^i\| \\ &\leq \sum_{i=1}^j \frac{\|\mathbf{d}_k^i\|^2 \|\nabla Q_1\|}{|\mathbf{d}_k^{iT} H_k \mathbf{d}_k^i|} \\ &= \sum_{i=1}^j \frac{\|\mathbf{d}_k^i\|^2}{|\mathbf{d}_k^{iT} H_k \mathbf{d}_k^i|} \|\nabla f(\mathbf{w}_k)\| \\ &\leq \frac{j}{M_2} \|\nabla f(\mathbf{w}_k)\| \\ &\leq \frac{n}{M_2} \|\nabla f(\mathbf{w}_k)\|, \end{aligned}$$

where the last inequality follows from the fact that Algorithm I always terminates in a number of steps that is bounded by n (see Theorem 5.1 of Nocedal and Wright (1999)). Thus, (xii) is proved to be satisfied with $a_2 = \frac{n}{M_2}$. \blacksquare

From Proposition III we have that \mathbf{s}_k is a descent direction and this means that Armijo line search procedure terminates finitely (see for instance Nocedal and Wright (1999) for details). Then we obtain global convergence by exploiting Armijo line search properties. In the following theorem we relax the assumption that \mathbf{s}_k is obtained by Algorithm I. In fact, the Theorem IV is still valid for any procedure that yields a direction \mathbf{s}_k that satisfies (xi) and (xii).

Theorem IV *Let $f \in LC^1$ and let \mathcal{L}_1 defined in (viii) be compact. Let $\{\mathbf{w}_k\}$ be a sequence generated by Algorithm II, where at each step k the direction in Step 3 is obtained by any procedure that yields a \mathbf{s}_k that satisfies (xi) and (xii). Then,*

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{w}_k)\| = 0. \quad (\text{xviii})$$

Proof From the Armijo line search instruction we get that

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \gamma \omega_k \nabla f(\mathbf{w}_k)^T \mathbf{s}_k, \quad (\text{xix})$$

and from this we obtain that

$$\begin{aligned}
f(\mathbf{w}_1) - f(\mathbf{w}_k) &= \sum_{i=1}^k f(\mathbf{w}_i) - f(\mathbf{w}_{i+1}) \\
&\geq -\gamma \sum_{i=1}^k \omega_i \nabla f(\mathbf{w}_i)^T \mathbf{s}_i \\
&\geq \gamma \sum_{i=1}^k \omega_i a_1 \|\nabla f(\mathbf{w}_i)\|^2,
\end{aligned} \tag{xx}$$

where the last inequality follows from (xi). By contradiction, there exists a subsequence $K \subset \{1, 2, \dots\}$ for which

$$\|\nabla f(\mathbf{w}_k)\| \geq \epsilon > 0, \quad \forall k \in K. \tag{xxi}$$

From Armijo line search properties we get that $\mathbf{w}_k \in \mathcal{L}_1$ and thus the left part of inequality (xx) is bounded from above, which means that also the right part of (xx) is bounded from above. Note that the right part of (xx) is also positive, thus,

$$\lim_{k \rightarrow \infty} \gamma \sum_{i=1}^k \omega_i a_1 \|\nabla f(\mathbf{w}_i)\| < \infty.$$

Since $\|\nabla f(\mathbf{w}_k)\|$ is bounded from below because of (xxi), we get that

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \sum_{i=1}^k \omega_i < \infty,$$

which, in turn, brings to

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \omega_k = 0. \tag{xxii}$$

Because of (xxii), it is not possible to have an infinite subsequence of K in which Armijo always terminates with $\omega_k = 1$. Thus, from Armijo line search properties there exists a \hat{k} for which $\forall k > \hat{k}, k \in K$ Armijo always terminates with a ω_k such that

$$f(\mathbf{w}_k + \frac{\omega_k}{\delta} \mathbf{s}_k) > f(\mathbf{w}_k) + \gamma \frac{\omega_k}{\delta} \nabla f(\mathbf{w}_k)^T \mathbf{s}_k.$$

From the Mean Value Theorem, applied on the above inequality, we obtain that

$$\nabla f(\mathbf{z}_k)^T \mathbf{s}_k > \gamma \nabla f(\mathbf{w}_k)^T \mathbf{s}_k, \tag{xxiii}$$

where $\mathbf{z}_k = \mathbf{w}_k + \theta_k \frac{\omega_k}{\delta} \mathbf{s}_k$ and $\theta_k \in [0, 1]$. Since $\{\mathbf{w}_k\}_K \in \mathcal{L}_1$, from K we can extract another infinite subsequence (redefined K) for which

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \mathbf{w}_k = \hat{\mathbf{w}}. \tag{xxiv}$$

Now, again from the fact that $\{\mathbf{w}_k\}_K \in \mathcal{L}_1$, we have that there exists a constant $M > 0$ such that $\|\nabla f(\mathbf{w}_k)\| < M \ \forall \mathbf{w}_k \in \mathcal{L}_1$ and thus, together with (xii) we get

$$\|\mathbf{s}_k\| \leq a_2 \|\nabla f(\mathbf{w}_k)\| \leq a_2 M. \quad (\text{xxv})$$

Thus, from (xxiv), (xxii) and (xxv) we get that

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \mathbf{z}_k = \lim_{\substack{k \rightarrow \infty \\ k \in K}} \mathbf{w}_k + \theta_k \frac{\omega_k \mathbf{s}_k}{\delta} = \hat{\mathbf{w}}. \quad (\text{xxvi})$$

From (xxv) we have that $\{\mathbf{s}_k\}_K$ is limited, so from K we can extract another infinite subsequence (redefined K) for which

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \mathbf{s}_k = \hat{\mathbf{s}}. \quad (\text{xxvii})$$

Moreover, from (xi) and (xxi) we have

$$\nabla f(\hat{\mathbf{w}})^T \hat{\mathbf{s}} \leq -a_1 \|\nabla f(\hat{\mathbf{w}})\|^2 \leq -a_1 \epsilon^2 < 0. \quad (\text{xxviii})$$

Finally, from (xxiii), (xxvi) and (xxvii) we get that

$$\gamma \nabla f(\hat{\mathbf{w}})^T \hat{\mathbf{s}} = \lim_{\substack{k \rightarrow \infty \\ k \in K}} \gamma \nabla f(\mathbf{w}_k)^T \mathbf{s}_k \leq \lim_{\substack{k \rightarrow \infty \\ k \in K}} \nabla f(\mathbf{z}_k)^T \mathbf{s}_k = \nabla f(\hat{\mathbf{w}})^T \hat{\mathbf{s}},$$

which is absurd, because from (xxviii) we have $\nabla f(\hat{\mathbf{w}})^T \hat{\mathbf{s}} < 0$ and $\gamma < 1$. ■

II. Details on Truncation Criteria in TNCG

We first recall all the *ratio*:

- residual:

$$ratio = \frac{\|\mathbf{g}_k + H_k \mathbf{s}_k^j\|}{\|\mathbf{g}_k\|}. \quad (\text{xxix})$$

- residual_{l1}:

$$ratio = \frac{\|\mathbf{g}_k + H_k \mathbf{s}_k^j\|_1}{\|\mathbf{g}_k\|_1}. \quad (\text{xxx})$$

- quadratic:

$$ratio = \frac{(Q_j - Q_{j-1})}{Q_j/j}, \quad (\text{xxxix})$$

Now we recall all the forcing sequences:

- constant:

$$\eta_k = c_0 \quad \text{with } c_0 \in (0, 1). \quad (\text{xxxixii})$$

- adaptive:

$$\eta_k = \min\{c_1; c_2 \|\mathbf{g}_k\|^{c_3}\}, \quad \text{with } c_1 \in (0, 1), c_2 > 0, c_3 \in (0, 1]. \quad (\text{xxxixiii})$$

- adaptive_{l1}:

$$\eta_k = \min\{c_1; c_2 \|\mathbf{g}_k\|_1^{c_3}\}, \quad \text{with } c_1 \in (0, 1), c_2 > 0, c_3 \in (0, 1]. \quad (\text{xxxixiv})$$

Algorithm III: PCG for solving (xxxv). Assume M has been factorized to EE^T .

```

1 Given  $\eta_k < 1, \Delta_k > 0$ , let  $\hat{\mathbf{s}} = \mathbf{0}, \hat{\mathbf{r}} = \hat{\mathbf{d}} = -E^{-1}\mathbf{g}_k, \gamma = \hat{\mathbf{r}}^T \hat{\mathbf{r}}$ 
2 while True do
3    $\hat{\mathbf{v}} \leftarrow E^{-1}(H_k(E^{-T}\hat{\mathbf{d}})),$ 
4    $\alpha \leftarrow \|\hat{\mathbf{r}}\|^2 / (\hat{\mathbf{d}}^T \hat{\mathbf{v}})$ 
5    $\hat{\mathbf{s}} \leftarrow \hat{\mathbf{s}} + \alpha \hat{\mathbf{d}}$ 
6    $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{r}} - \alpha \hat{\mathbf{v}}$ 
7   if  $\|\hat{\mathbf{r}}\| < \eta_k \|\hat{\mathbf{g}}_k\|$  then
8     return  $\mathbf{s}_k = E^{-T}\hat{\mathbf{s}}$ 
9    $\gamma^{\text{new}} \leftarrow \hat{\mathbf{r}}^T \hat{\mathbf{r}}$ 
10   $\beta \leftarrow \gamma^{\text{new}} / \gamma$ 
11   $\hat{\mathbf{d}} \leftarrow \hat{\mathbf{r}} + \beta \hat{\mathbf{d}}$ 
12   $\gamma \leftarrow \gamma^{\text{new}}$ 

```

III. Preconditioning

In this section we will show that Algorithm 2 from the original paper is equivalent to applying Algorithm III and then obtaining \mathbf{s} from $\mathbf{s} = E^{-1}\hat{\mathbf{s}}$. For sake of simplicity in this section and in Algorithm III we will not use the iteration counter j and k will only be reported on H_k and \mathbf{g}_k . Let us write the instructions of the original CG method (Algorithm I) applied to the system

$$E^{-1}H_kE^{-T}\hat{\mathbf{s}} = -E^{-1}\mathbf{g}_k. \quad (\text{xxxv})$$

At the internal generic CG step we have

$$\hat{\mathbf{s}} = \hat{\mathbf{s}} + \alpha \hat{\mathbf{d}}, \quad (\text{xxxvi})$$

where

$$\alpha = \frac{\|\hat{\mathbf{r}}\|^2}{\hat{\mathbf{d}}^T E^{-1}H_kE^{-T}\hat{\mathbf{d}}}$$

with

$$\hat{\mathbf{r}} = -E^{-1}H_kE^{-T}\hat{\mathbf{s}} - E^{-1}\mathbf{g}_k = E^{-1}(-H_kE^{-T}\hat{\mathbf{s}} - \mathbf{g}_k). \quad (\text{xxxvii})$$

Pre-multiplying (xxxvi) by E^{-1} , using the transformation $\mathbf{s} = E^{-1}\hat{\mathbf{s}}$ and the definition

$$\mathbf{d} := E^{-1}\hat{\mathbf{d}} \quad (\text{xxxviii})$$

we get

$$\mathbf{s} = \mathbf{s} + \alpha \mathbf{d}.$$

Again using the transformation $\mathbf{s} = E^{-1}\hat{\mathbf{s}}$, (xxxvii) can be re-written to get

$$\hat{\mathbf{r}} = E^{-1}(-H_k\mathbf{s} - \mathbf{g}_k) = E^{-1}\mathbf{r}, \quad (\text{xxxix})$$

which together with (xxxviii) brings to

$$\alpha = \frac{\mathbf{r}^T M^{-1} \mathbf{r}}{\mathbf{d}^T \mathbf{v}}, \quad \text{where } \mathbf{v} = H_k \mathbf{d}.$$

Moreover we have that

$$\hat{\mathbf{d}} = \hat{\mathbf{r}} + \beta \hat{\mathbf{d}}, \quad (\text{xl})$$

where

$$\beta = \frac{\hat{\mathbf{r}}_{\text{new}}^T \hat{\mathbf{r}}_{\text{new}}}{\hat{\mathbf{r}}^T \hat{\mathbf{r}}} = \frac{\mathbf{r}_{\text{new}}^T M^{-1} \mathbf{r}_{\text{new}}}{\mathbf{r}^T M^{-1} \mathbf{r}}.$$

Pre-multiplying (xl) by E^{-1} , from (xxxix) and (xxxviii) we get

$$\mathbf{d} = \mathbf{r} + \beta \mathbf{d}.$$

Finally, defining the vector

$$\mathbf{z} := M^{-1} \mathbf{r} \quad \text{and the scalar } \gamma := \mathbf{r}^T \mathbf{z},$$

α and β can be re-written as

$$\alpha = \frac{\mathbf{r}^T \mathbf{z}}{\mathbf{d}^T \mathbf{v}} \quad \text{and} \quad \beta = \frac{\gamma^{\text{new}}}{\gamma}.$$

IV. Complete Numerical Results

For experiments we use two-class classification problems shown in Table I. Except for `yahookr` and `yahoojp`, others are available from LIBSVM Data Sets (2007). Table I shows the statistics of each data set used in our experiments. For a fair evaluation, all different

Data sets	#instances	#features	density	$\log_2(C_{\text{Best}})$	
				LR	L2
<code>news20</code>	19,996	1,355,191	0.0336%	9	3
<code>w8a</code>	49,749	300	3.8834%	8	2
<code>covtype</code>	581,012	54	TODO	-23	-26
<code>rcv1</code>	20,242	47,236	0.1568%	6	0
<code>url</code>	2,396,130	3,231,962	0.0036%	-7	-10
<code>real-sim</code>	72,309	20,958	0.2448%	3	-1
<code>yahoojp</code>	176,203	832,026	0.0160%	3	-1
<code>yahookr</code>	460,554	3,052,939	0.0111%	6	1
<code>kdd2010a</code>	8,407,752	20,216,831	0.0001%	-3	-5
<code>kdd2010b</code>	19,264,097	29,890,095	0.0001%	-1	-4
<code>HIGGS</code>	11,000,000	28	92.1057%	-6	-12
<code>webspam</code>	350,000	16,609,143	0.0220%	2	-3
<code>criteo</code>	45,840,617	1,000,000	0.0039%	-15	-12
<code>kdd2012</code>	149,639,105	54,686,452	0.00002%	-4	-11

Table I: Data statistics. The density is the average number of non-zeros per instance. C_{Best} is the regularization parameter selected by cross validation.

settings are implemented based on the LIBLINEAR package¹. To simulate the practical

1. Exceptions are for Scikit and ScikitArmijo of Figure 1 from the original paper, where Scikit-learn is used. In particular, to train a `LogisticRegression` model, the user should select the `newton-cg` option for the parameter `solver`. By default Scikit-learn solves a slightly different optimization problem with a bias term. To solve (i), the bias term is thus ignored by setting the option `fit_intercept` to `False`.

use we conduct a five-fold cross-validation to select the regularization parameter C_{Best} that achieves the best validation accuracy (see Table I for the resulting C). Then in experiments we consider $C = C_{\text{Best}} \times \{1, 100\}$.

In figures below we show the total number of CG steps needed to solve the training problem (i). On the y axis we report the relative reduction of the function value, computed by

$$\frac{f(\mathbf{w}_k) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)},$$

where \mathbf{w}^* is the optimal solution² of (i). The four horizontal lines in figures below indicate places where the following stopping condition is met respectively with $\epsilon = \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$

$$\|\mathbf{g}_k\| \leq \epsilon \frac{\min\{\#\text{pos}, \#\text{neg}\}}{l} \cdot \|\mathbf{g}_1\|, \quad (\text{xli})$$

where $\#\text{pos}$, $\#\text{neg}$ are the numbers of positive- and negative-labeled instances respectively, and $\epsilon > 0$ is the thresholding constant that controls the precision of the training procedure. Note that (xli) with $\epsilon = 10^{-2}$ is the outer stopping condition employed in LIBLINEAR. Thus the behavior before 10^{-1} and after 10^{-4} is not crucial, since the training would be stopped too early or too late. For the Armijo line search, we follow Hsia et al. (2017) to have $\delta = 0.5$ for deciding the step size $\alpha \in \{1, \delta, \delta^2, \dots\}$ and $\gamma = 0.01$ in the sufficient decrease condition (vi).

IV.1 Selection of the Constant Threshold in the Inner Stopping Condition

See complete results in Figures i and ii.

IV.2 Comparison Between Adaptive and Constant Forcing Sequences

See complete results in Figures iii and iv.

IV.3 Robustness of Rules with Adaptive Forcing Sequences and Comparison with Trust-Region

See complete results in Figures v and vi. On the dataset `covtype` with $C = 100C_{\text{Best}}$ both the setting `quadada` and `quadada_ll` are facing slow convergence issues because of the threshold c_1 . Comparing with `resada` it seems that in early iterations the Newton equation needed to be solved more accurately. For this reason we investigated the choice of the threshold c_1 for `quadada` in Section IV.7.

IV.4 Analysis via Conditions for Global Convergence

See complete results in Figures vii and viii.

2. In practice \mathbf{w}^* is not available, so we run enough iterations to get a good approximate solution.

IV.5 Preconditioning: Comparison Against Rules with Constant Forcing Sequences

See complete results in Figures ix and ix.

IV.6 Preconditioning: Comparison of Rules with Adaptive Forcing Sequences and Comparison with Trust Region

See complete results in Figures xi and xii.

IV.7 Additional Results in the Preconditioned Case: Selection of the c_1 Threshold in the Quadratic Adaptive Rule

In this section, results in Figures xiii and xiv are obtained by comparing

- quadada01: $ratio = (xxxix)$, $\eta_k = (xxxiii)$, $c_1 = 0.1$, $c_2 = 1$, $c_3 = 0.5$;
- quadada01_p: $ratio = (xxxix)$, $\eta_k = (xxxiii)$, $c_1 = 0.1$, $c_2 = 1$, $c_3 = 0.5$, preconditioned;
- quadada: $ratio = (xxxix)$, $\eta_k = (xxxiii)$, $c_1 = 0.5$, $c_2 = 1$, $c_3 = 0.5$;
- quadada_p: $ratio = (xxxix)$, $\eta_k = (xxxiii)$, $c_1 = 0.5$, $c_2 = 1$, $c_3 = 0.5$ preconditioned.

From Figures xiii and xiv we can make the following observations:

- Using a smaller c_1 threshold ($c_1 = 0.1$) is solving the issues encountered by quadada on `covtype` with $C = 100C_{Best}$ in Figure vi.
- Both the settings quadada01 and quadada01_p are generally slower than quadada and quadada_p.
- The slow convergence issue encountered by quadada on `covtype` with $C = 100C_{Best}$ in Figure vi is also solved by the use of preconditioning (quadada_p).

IV.8 Additional Results in the Preconditioned Case: Comparison Against a Residual Constant Rule with Higher Constant Threshold

In this section, results in Figures xv and xvi are obtained by comparing

- rescons09: $ratio = (xxix)$, $\eta_k = (xxxii)$, $c_0 = 0.9$;
- rescons09_p: $ratio = (xxix)$, $\eta_k = (xxxii)$, $c_0 = 0.9$, preconditioned;
- quadada: $ratio = (xxxix)$, $\eta_k = (xxxiii)$;
- quadada_p: $ratio = (xxxix)$, $\eta_k = (xxxiii)$, preconditioned.

From Figures xv and xvi we can make the following observations:

- The settings rescons09_p and quadada_p are often obtaining very similar speed of convergence.

- There are difference on $C = C_{\text{Best}}$ in which `rescons09_p` is faster than `quadada_p` (`rcv1`, `real-sim`) and `quadada_p` is faster than `rescons09_p` (`covtype` and `criteo` on the fourth line), but the detachment is less than 20 CG steps.
- By looking at `yahookr` with $C = C_{\text{Best}}$, we can see that `rescons09_p` is noticeably faster than `quadada_p` (at the third line is around twice faster, by 150 CG steps). The opposite situation is instead happening on `yahookr` with $C = 100C_{\text{Best}}$, where being twice faster on the third line means saving around 750 CG steps.

IV.9 Additional Results in the Preconditioned Case: Robustness of Rules with Adaptive Forcing Sequences

In this section, results in Figures xvii and xviii are obtained by comparing

- `resada_l1`: $ratio = (\text{xxx})$, $\eta_k = (\text{xxxiv})$;
- `rescons_l1_p`: $ratio = (\text{xxx})$, (xxxiv) , preconditioned;
- `resada`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxiii})$;
- `resada_p`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxiii})$, preconditioned.

From Figures xvii and xviii we can notice that `rescons_l1_p` is often slower than `resada_p`, so also in the preconditioned case there is no reason to employ L1-norm instead of L2-norm.

IV.10 Results for L2-loss SVM

In this section, results in Figures xix and xx are obtained by employing the ξ_{L2} loss as defined in (ii) and we compare:

- `tr_rescons`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxii})$, $c_0 = 0.1$, Trust Region instead of line search;
- `rescons09`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxii})$, $c_0 = 0.9$;
- `resada`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxiii})$;
- `quadada`: $ratio = (\text{xxx})$, $\eta_k = (\text{xxxiii})$.

IV.11 Results for L2-loss SVM in the Preconditioned Case

In this section, results in Figures xxi and xxii are obtained by employing the ξ_{L2} loss as defined in (ii) and we compare:

- `tr_rescons_p`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxii})$, $c_0 = 0.1$, Trust Region instead of line search, preconditioned;
- `rescons09_p`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxii})$, $c_0 = 0.9$, preconditioned;
- `resada_p`: $ratio = (\text{xxix})$, $\eta_k = (\text{xxxiii})$, preconditioned;
- `quadada_p`: $ratio = (\text{xxx})$, $\eta_k = (\text{xxxiii})$, preconditioned.

Appendix A. Gradient and Hessian Details of the Losses from (ii)

Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_l]^T$ be the data matrix.

A.1 Logistic Loss

In this section we compute the gradient and the Hessian of f when the loss employed is ξ_{LR} .

Let $\sigma(t) := (1 + e^{-t})^{-1}$.

The gradient of f is

$$\nabla f(\mathbf{w}) = \mathbf{w} + C \sum_{i=1}^l (\sigma(y_i \mathbf{w}^T \mathbf{x}_i) - 1) y_i \mathbf{x}_i$$

The Hessian of f is

$$\nabla^2 f(\mathbf{w}) = I + CX^T DX$$

with D a diagonal matrix such that $D_{ii}(\mathbf{w}) = \sigma(y_i \mathbf{w}^T \mathbf{x}_i)(1 - \sigma(y_i \mathbf{w}^T \mathbf{x}_i))$.

A.2 L2-loss SVM

In this section we compute the gradient and the Hessian of f when the loss employed is ξ_{L2} .

Let $\mathbf{y} = [y_1, \dots, y_l]^T$. The gradient of f is

$$\nabla f(\mathbf{w}) = \mathbf{w} + 2CX_{\mathcal{I}}{:}^T (X_{\mathcal{I}}{:} \mathbf{w} - \mathbf{y}_{\mathcal{I}}),$$

where $\mathcal{I} = \{i \in \{1, \dots, l\} | 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0\}$.

Now we define

$$\phi(t) = \max(0, 1 - t)^2, \quad \phi'(t) = \begin{cases} -2(1 - t) & \text{if } 1 - t > 0, \\ 0 & \text{otherwise,} \end{cases}$$

$$\partial_B \phi'(t) = \begin{cases} 2 & \text{if } 1 - t > 0, \\ 0 & \text{if } 1 - t < 0, \\ \{0, 2\} & \text{if } 1 - t = 0, \end{cases}$$

$$\partial \phi'(t) = \text{conv}\{\partial_B \phi'(t)\} = \begin{cases} 2 & \text{if } 1 - t > 0, \\ 0 & \text{if } 1 - t < 0, \\ \mu \text{ with } \mu \in [0, 2] & \text{if } 1 - t = 0, \end{cases}$$

where $\text{conv}\{\cdot\}$ is the convex hull of a set. Therefore, the generalized Hessian of f is

$$\partial \nabla f = I + CX^T DX, \quad \text{with } D \text{ a diagonal matrix such that } D_{ii}(\mathbf{w}) = \partial \phi'(y_i \mathbf{w}^T \mathbf{x}_i).$$

An easy way to calculate $\partial \nabla f$ in practice is to replace the original $D_{ii}(\mathbf{w})$ with the following

$$\bar{D}_{ii}(\mathbf{w}) = \begin{cases} 2, & \text{if } 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0, \\ 0, & \text{otherwise.} \end{cases}$$

A.3 Cost of Armijo, Wolfe and Exact Line Searches

Here we recall the cost of Armijo (Section 2.1 of Hsia et al. (2017)) and exact line search (Section 3.3 of Hsia et al. (2017)), and we show how to obtain the same save of operations even for Wolfe line search.

- Exact requires $\mathcal{O}(l)$ operations for each trial ω ,
- Armijo requires $\mathcal{O}(l)$ operations for each trial ω ,
- Wolfe requires $\mathcal{O}(l)$ operations for each trial ω .

To show the trick in computing the Wolfe condition, we first recall the same trick for Armijo (vi). For computing f at each ω , we can store each $\mathbf{x}_i^T \mathbf{w}_k$, $\mathbf{x}_i^T \mathbf{s}_k$, $\mathbf{w}_k^T \mathbf{s}_k$, $\mathbf{w}_k^T \mathbf{w}_k$ and $\mathbf{s}_k^T \mathbf{s}_k$ respectively in these scalars $x_i \bar{w}$, $x_i \bar{s}$, $\bar{w} \bar{s}$, $\bar{w} \bar{w}$ and $\bar{s} \bar{s}$ and compute f at $(\mathbf{w}_k + \omega \mathbf{s}_k)$ in the following way

$$\begin{aligned} f(\mathbf{w}_k + \omega \mathbf{s}_k) &= \frac{1}{2} (\mathbf{w}_k + \omega \mathbf{s}_k)^T (\mathbf{w}_k + \omega \mathbf{s}_k) + C \sum_{i=1}^l \xi(y_i (\mathbf{w}_k + \omega \mathbf{s}_k)^T \mathbf{x}_i) \\ &= \frac{1}{2} (\bar{w} \bar{w} + 2\omega \bar{w} \bar{s} + \omega^2 \bar{s} \bar{s}) + C \sum_{i=1}^l \xi(y_i (x_i \bar{w} + \omega x_i \bar{s})). \end{aligned}$$

Thus the cost of Armijo line search is $\mathcal{O}(l) \times$ (line search steps). The cost of Wolfe condition (vii) is also $\mathcal{O}(l) \times$ (line search steps), because we don't have to compute $\nabla f(\mathbf{w}_k + \omega \mathbf{s}_k)$ at each ω , but we can directly compute $\nabla f(\mathbf{w}_k + \omega \mathbf{s}_k)^T \mathbf{s}_k$. Thus, again by storing everything as above, we have:

$$\begin{aligned} \nabla f(\mathbf{w}_k + \omega \mathbf{s}_k)^T \mathbf{s}_k &= \mathbf{w}_k^T \mathbf{s}_k + \omega \mathbf{s}_k^T \mathbf{s}_k + C \sum_{i=1}^l \xi'(y_i (\mathbf{w}_k + \omega \mathbf{s}_k)^T \mathbf{x}_i) \cdot y_i \cdot \mathbf{x}_i^T \mathbf{s}_k \\ &= \bar{w} \bar{s} + \omega \bar{s} \bar{s} + C \sum_{i=1}^l \xi'(y_i (x_i \bar{w} + \omega x_i \bar{s})) \cdot y_i \cdot x_i \bar{s}. \end{aligned}$$

References

- Magnus Rudolph Hestenes and Eduard Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of Research of the National Bureau of Standards*, 49(1):409–436, 1952.
- Chih-Yang Hsia, Ya Zhu, and Chih-Jen Lin. A study on trust region update rules in Newton methods for large-scale linear classification. In *Proceedings of the Asian Conference on Machine Learning (ACML)*, 2017. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/newtron/newtron.pdf>.
- Ching-Pei Lee, Po-Wei Wang, Weizhu Chen, and Chih-Jen Lin. Limited-memory common-directions method for distributed optimization and its application on empirical risk minimization. In *Proceedings of SIAM International Conference on Data Mining (SDM)*, 2017. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/l-commdir/l-commdir.pdf>.

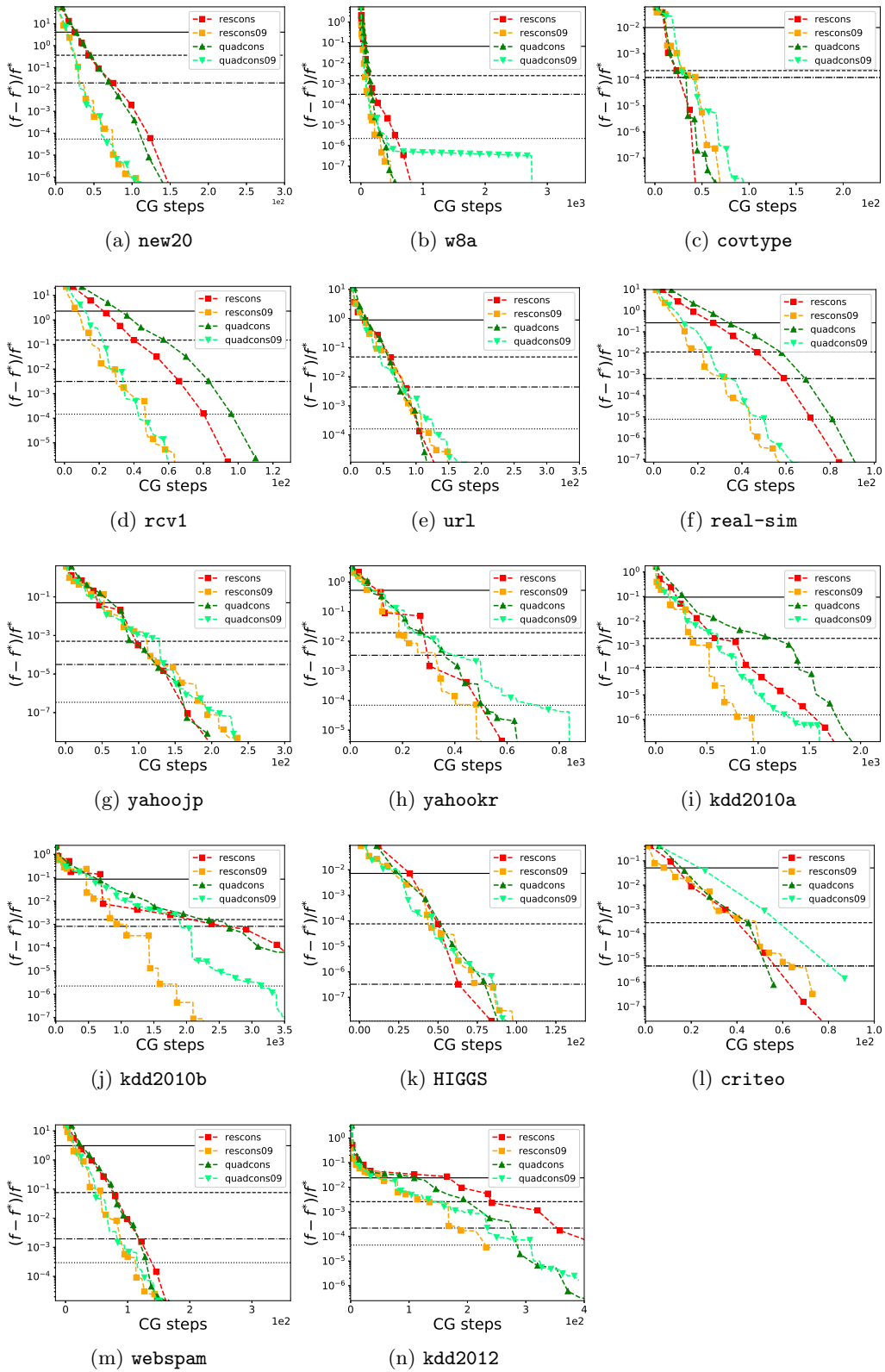


Figure i: A comparison of different constant thresholds in the inner stopping condition. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

LIBSVM Data Sets, 2007. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>.

Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer-Verlag, New York, NY, 1999.

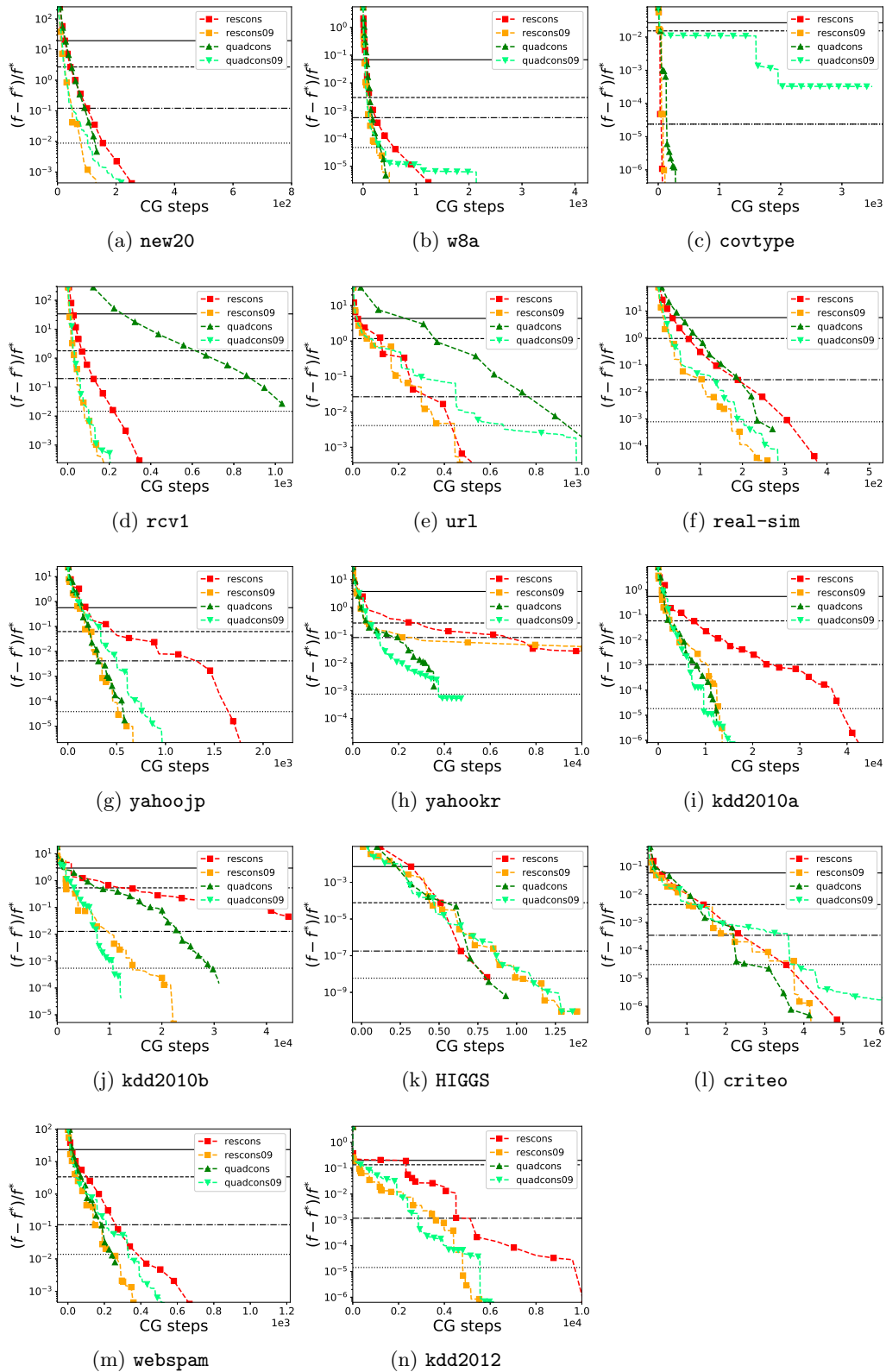


Figure ii: A comparison of different constant thresholds in the inner stopping condition. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

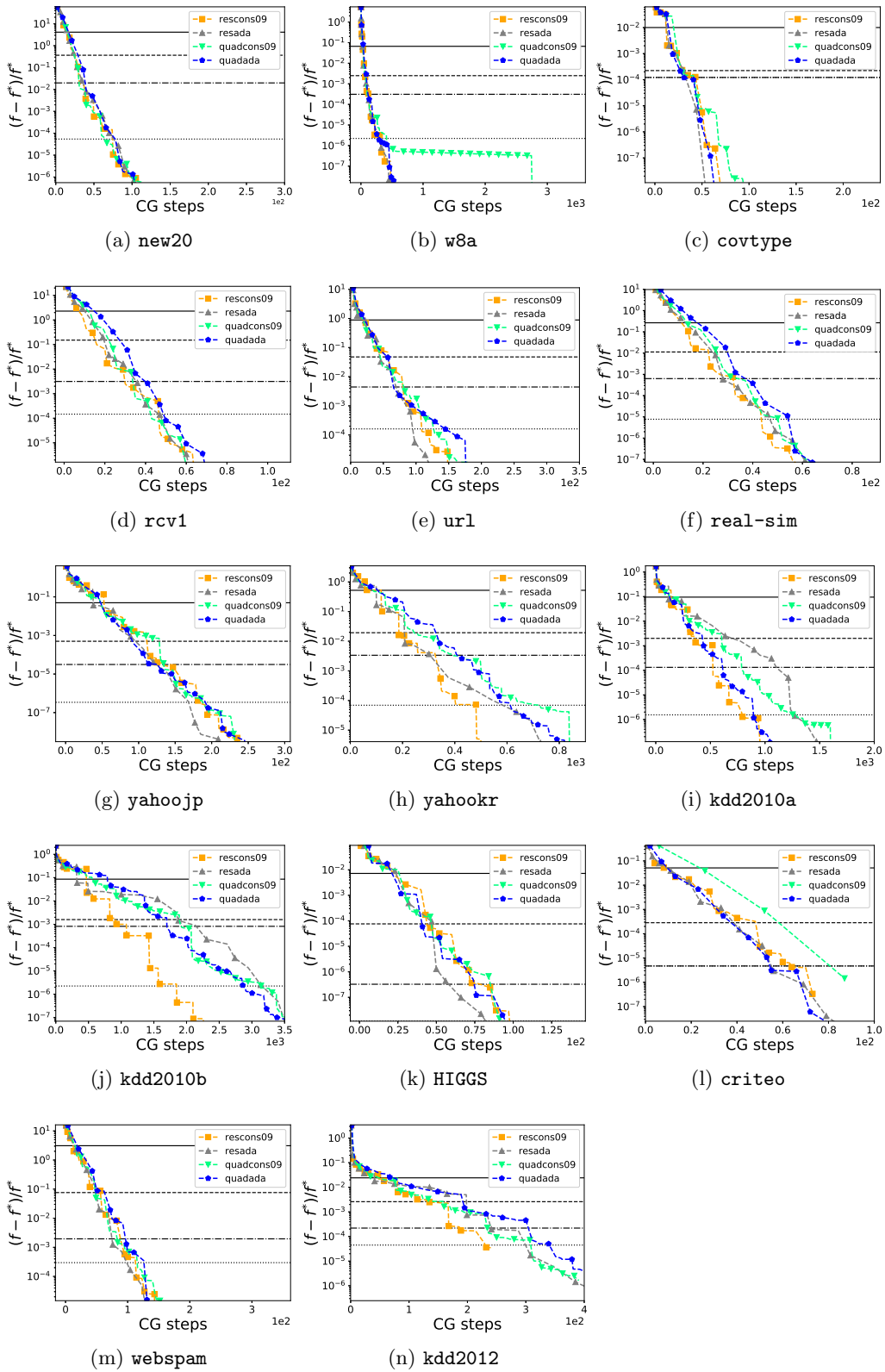


Figure iii: A comparison between adaptive and constant forcing sequences. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

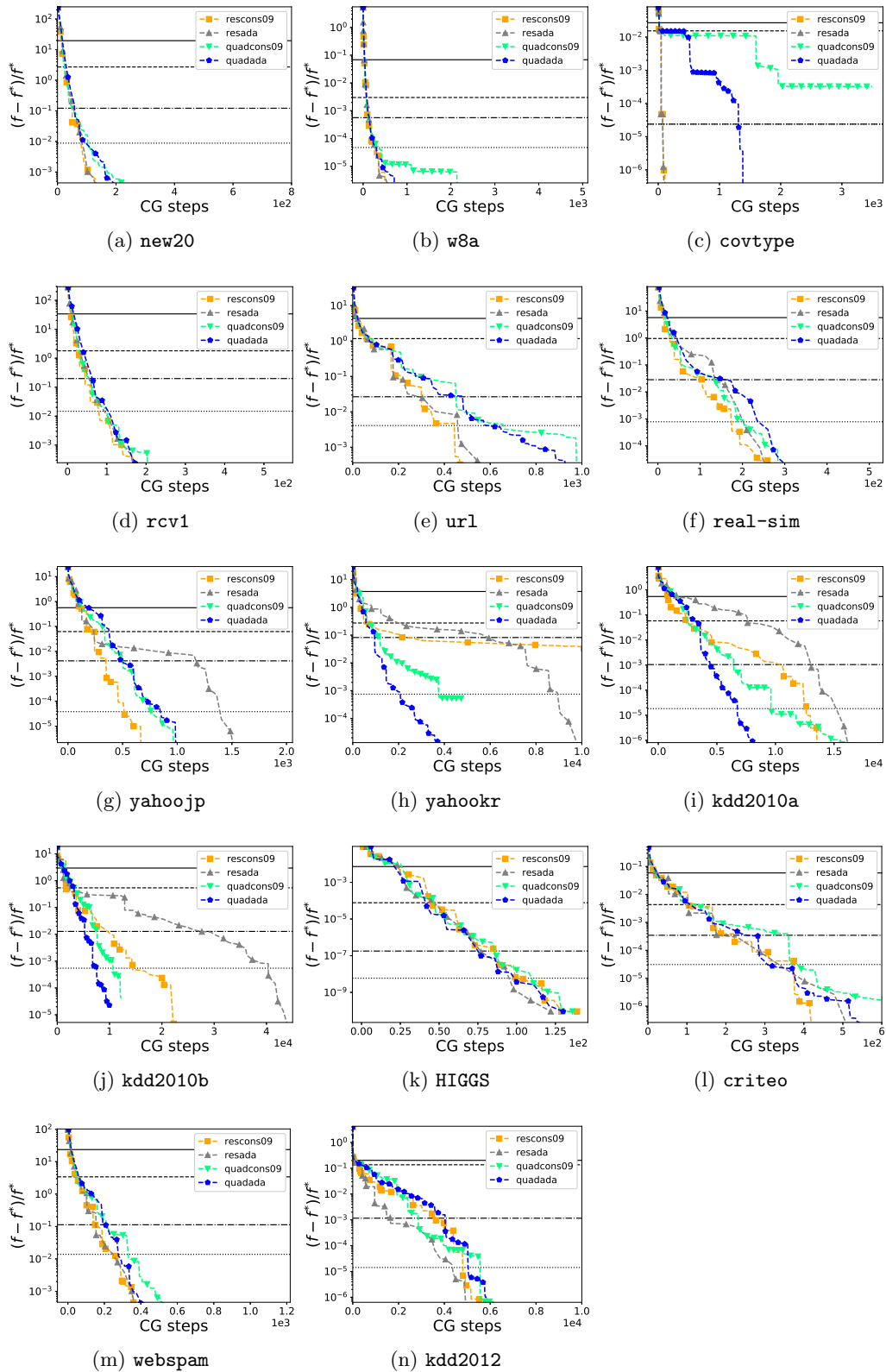


Figure iv: A comparison between adaptive and constant forcing sequences. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

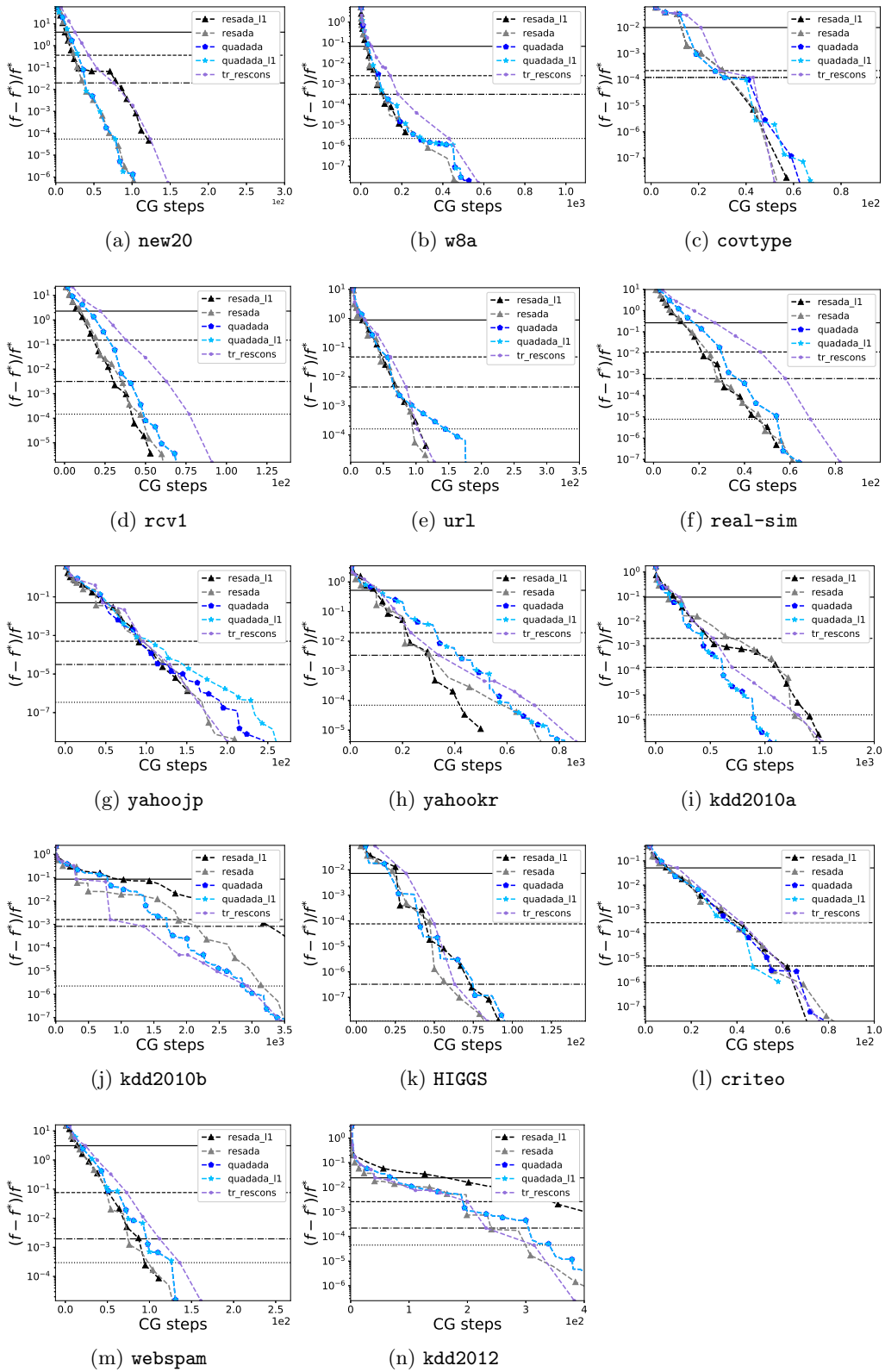


Figure v: An investigation on the robustness of adaptive rules and a comparison with the trust region approach. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

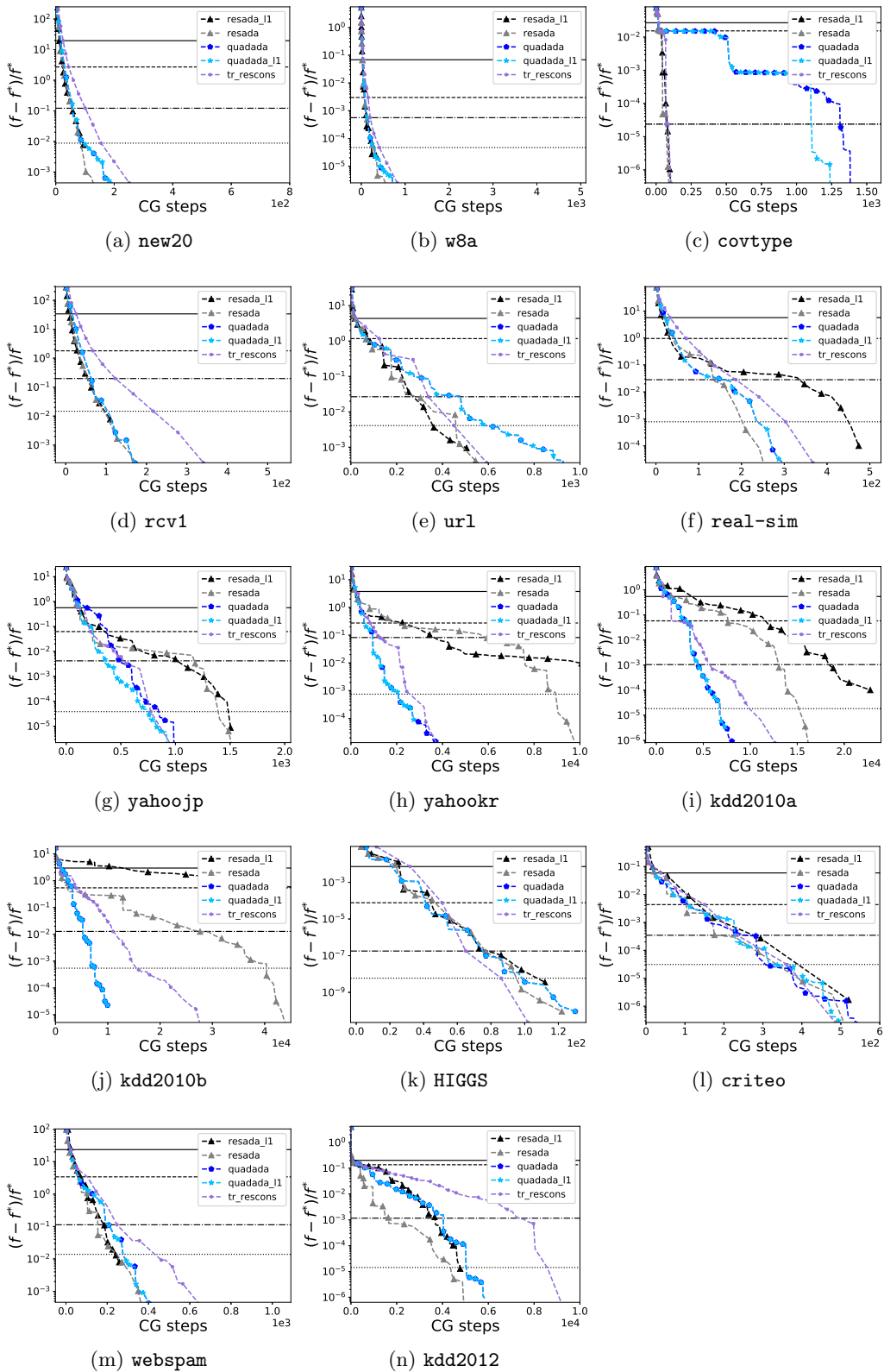


Figure vi: An investigation on the robustness of adaptive rules and a comparison with the trust region approach. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

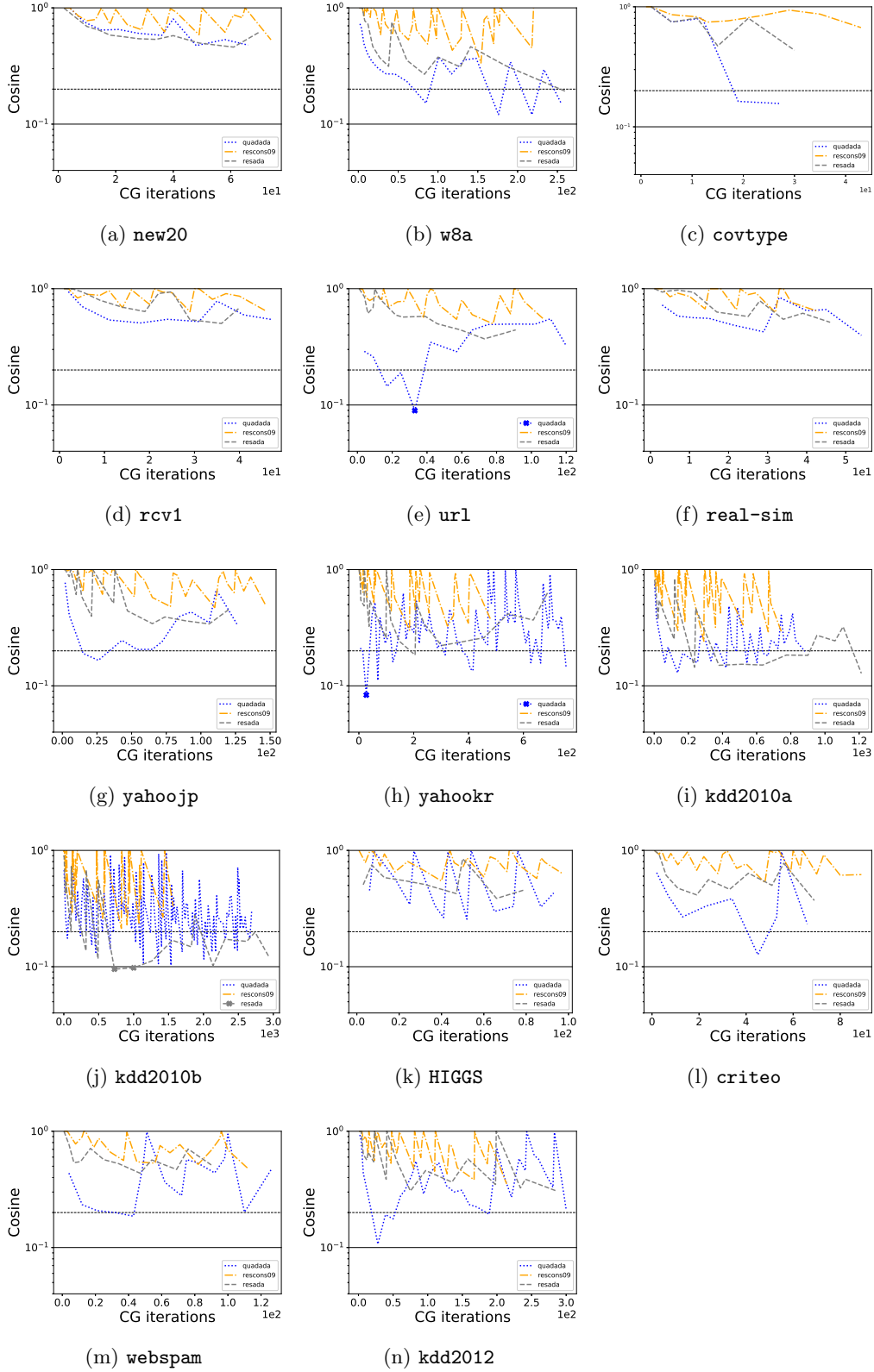


Figure VII: Cosines between the anti-gradient $-\nabla f(\mathbf{w}_k)$ and the resulting direction \mathbf{s}_k , using different truncation rules. Steps in which the cosine is below the 0.1 threshold are marked with a \times . The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

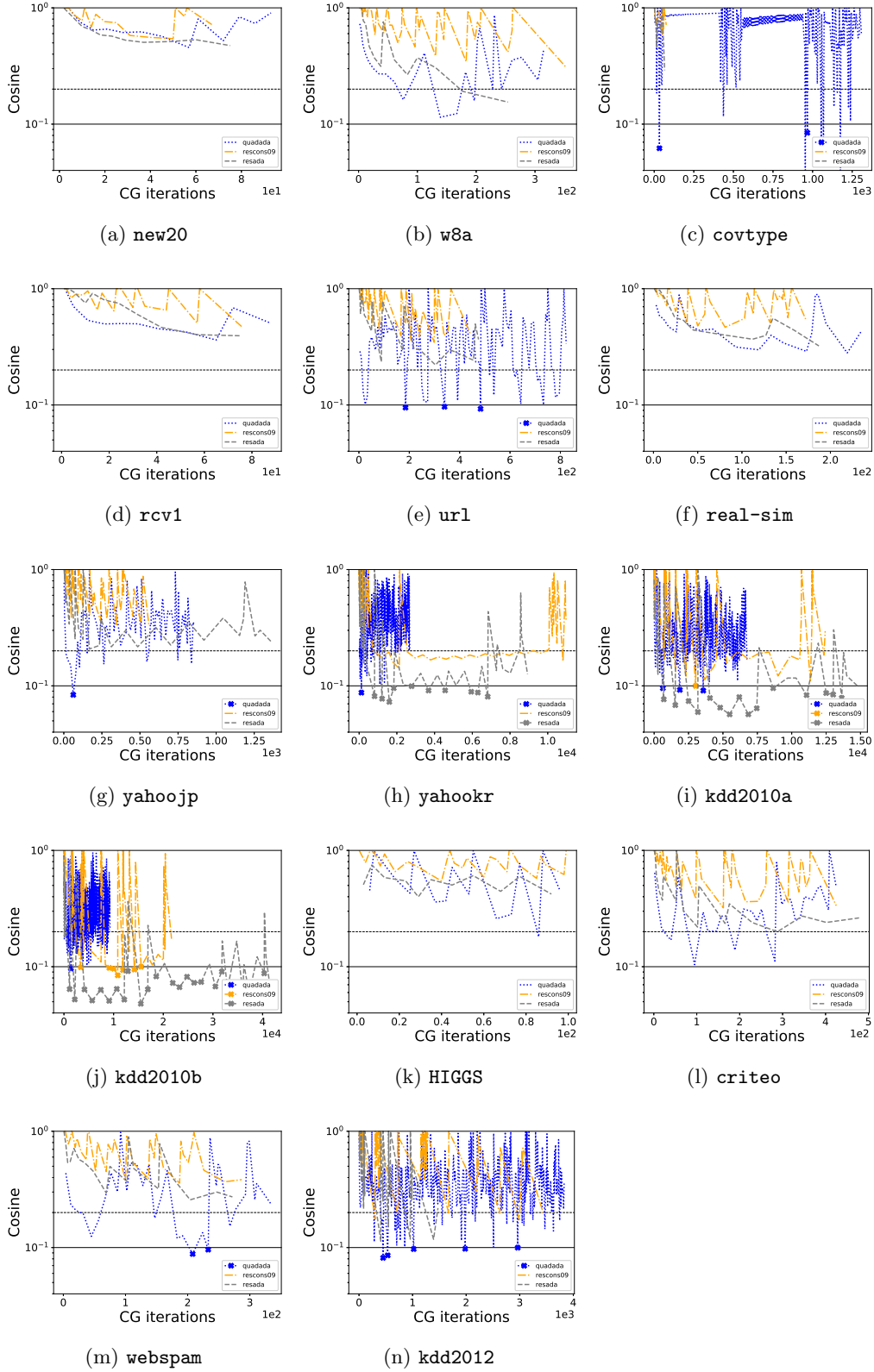


Figure VIII: Cosines between the anti-gradient $-\nabla f(\mathbf{w}_k)$ and the resulting direction \mathbf{s}_k , using different truncation rules. Steps in which the cosine is below the 0.1 threshold are marked with a \times . The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

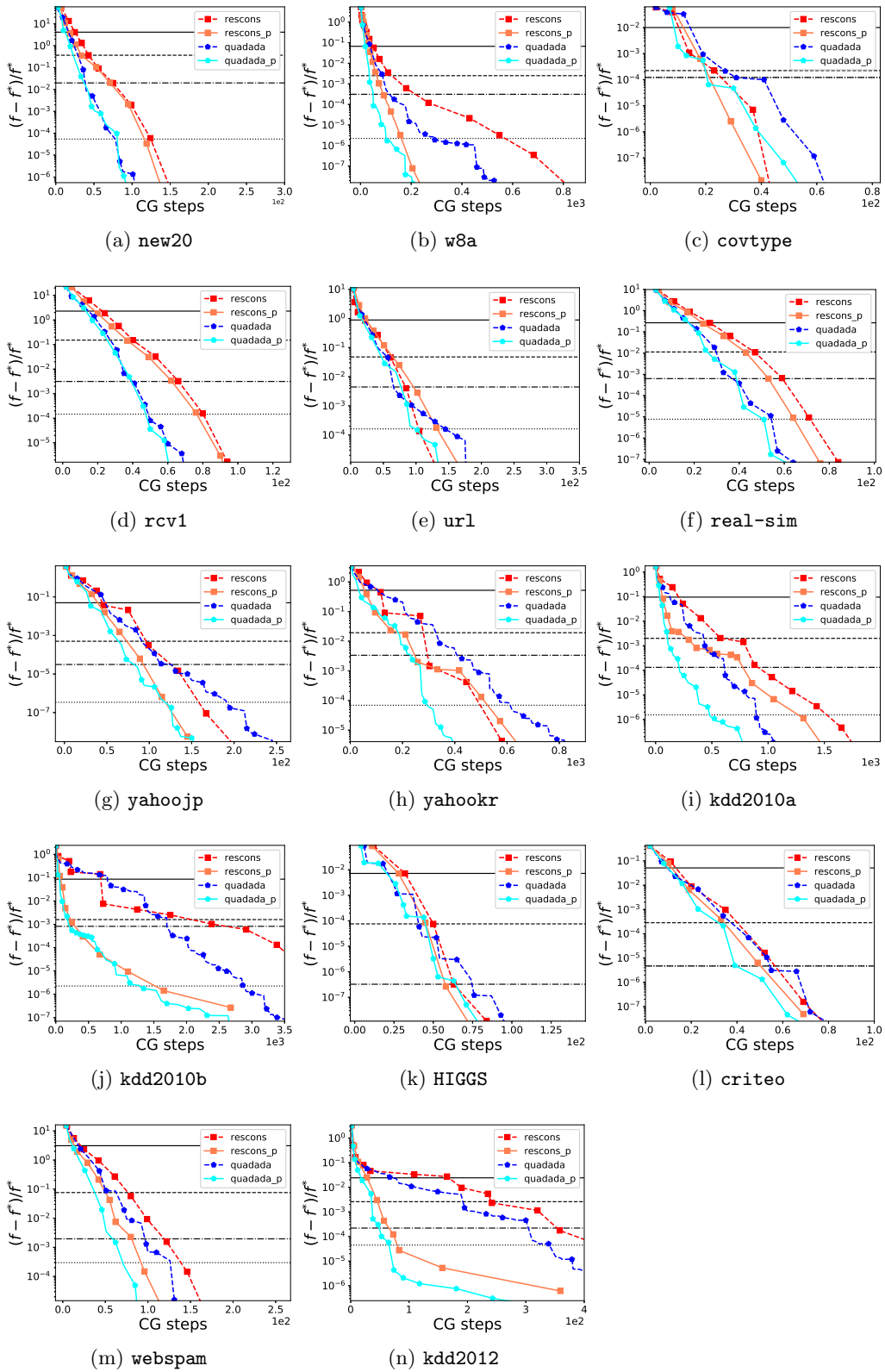


Figure ix: A comparison between adaptive and constant forcing sequences in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure. Loss=LR and $C = C_{\text{Best}}$.

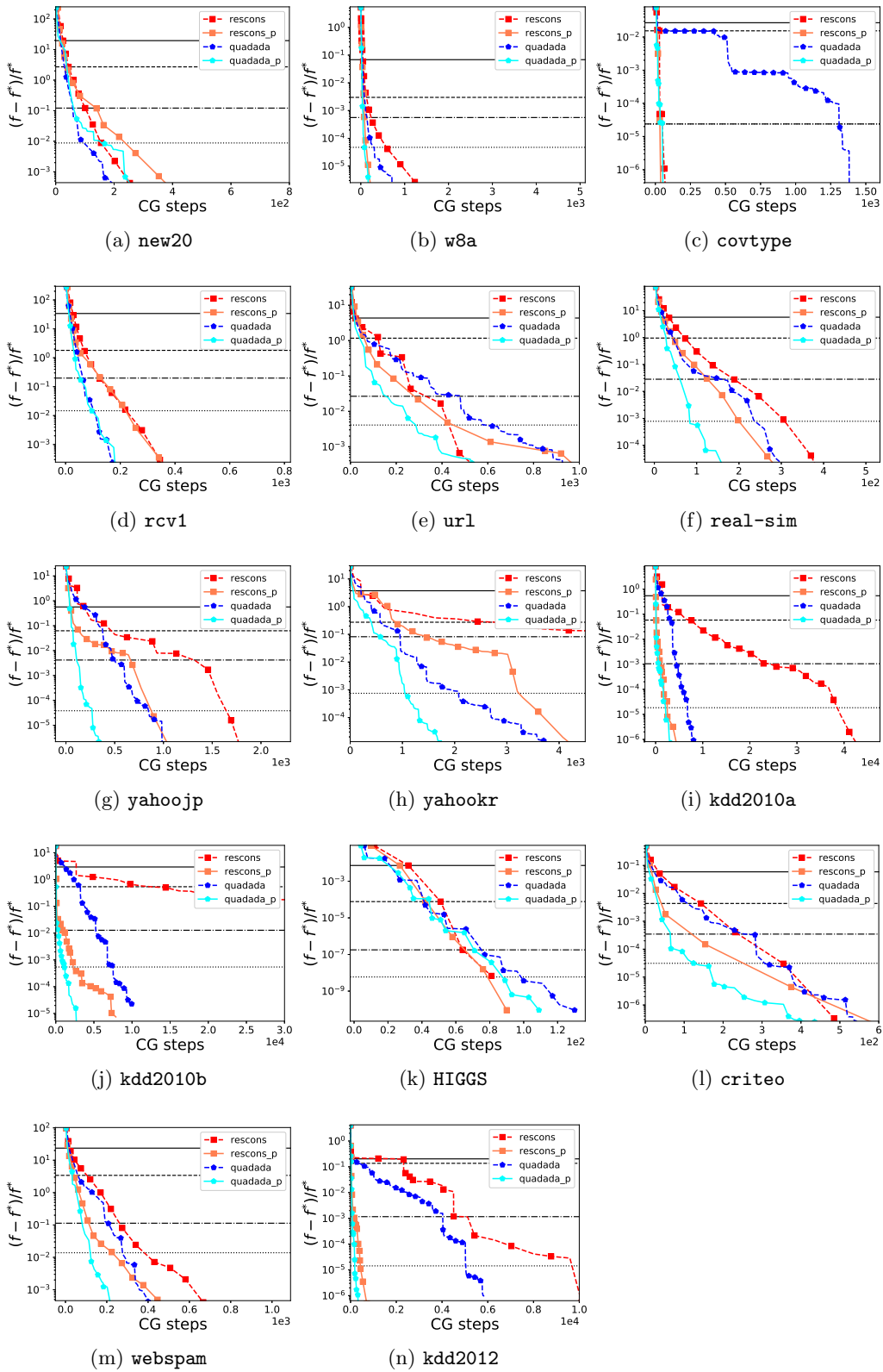


Figure x: A comparison between adaptive and constant forcing sequences in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

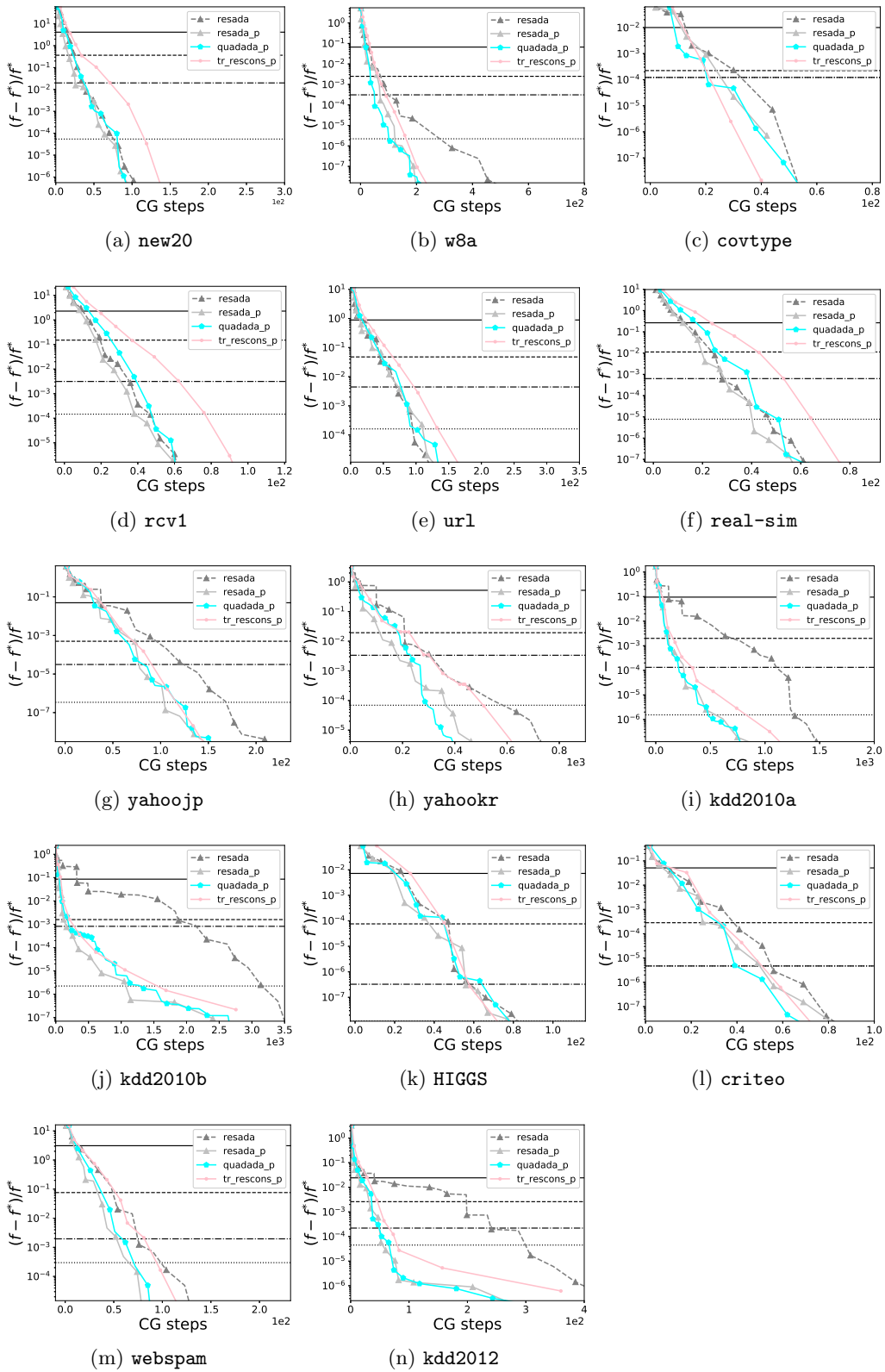


Figure xi: A comparison between adaptive rules and the trust region approach in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

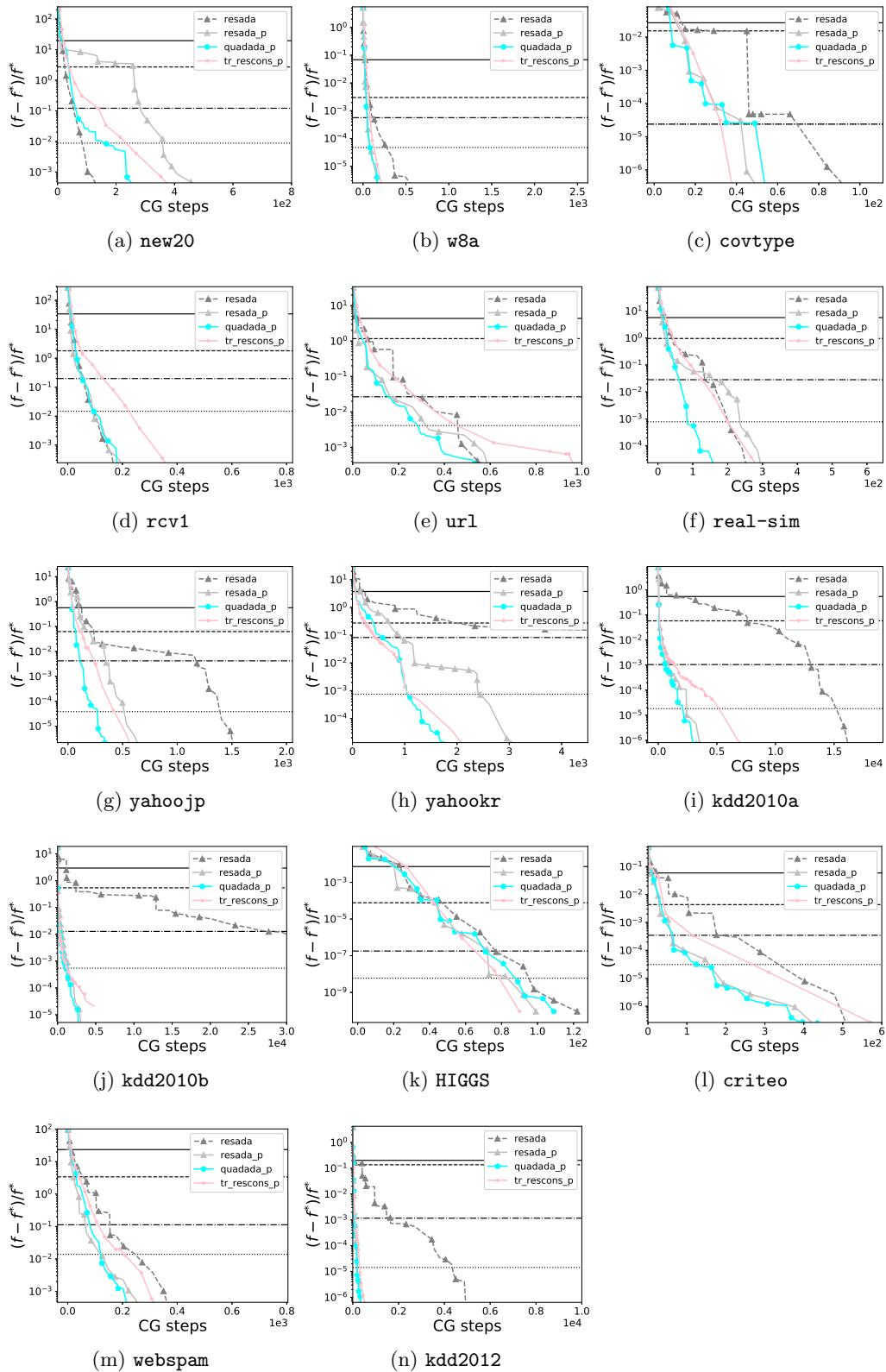


Figure xii: A comparison between adaptive rules and the trust region approach in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

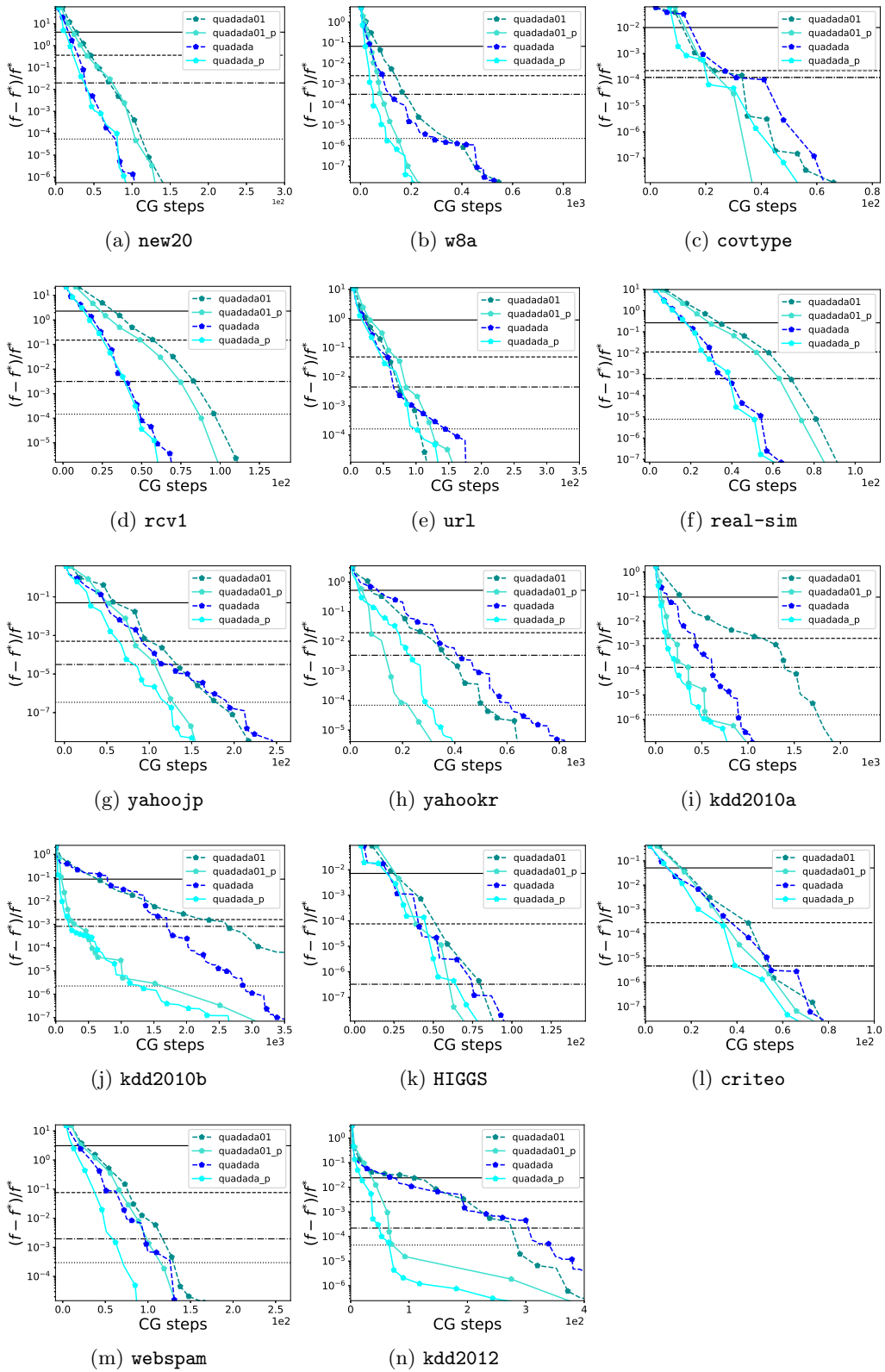


Figure xiii: A comparison of different early thresholds in the quadratic adaptive inner stopping condition. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

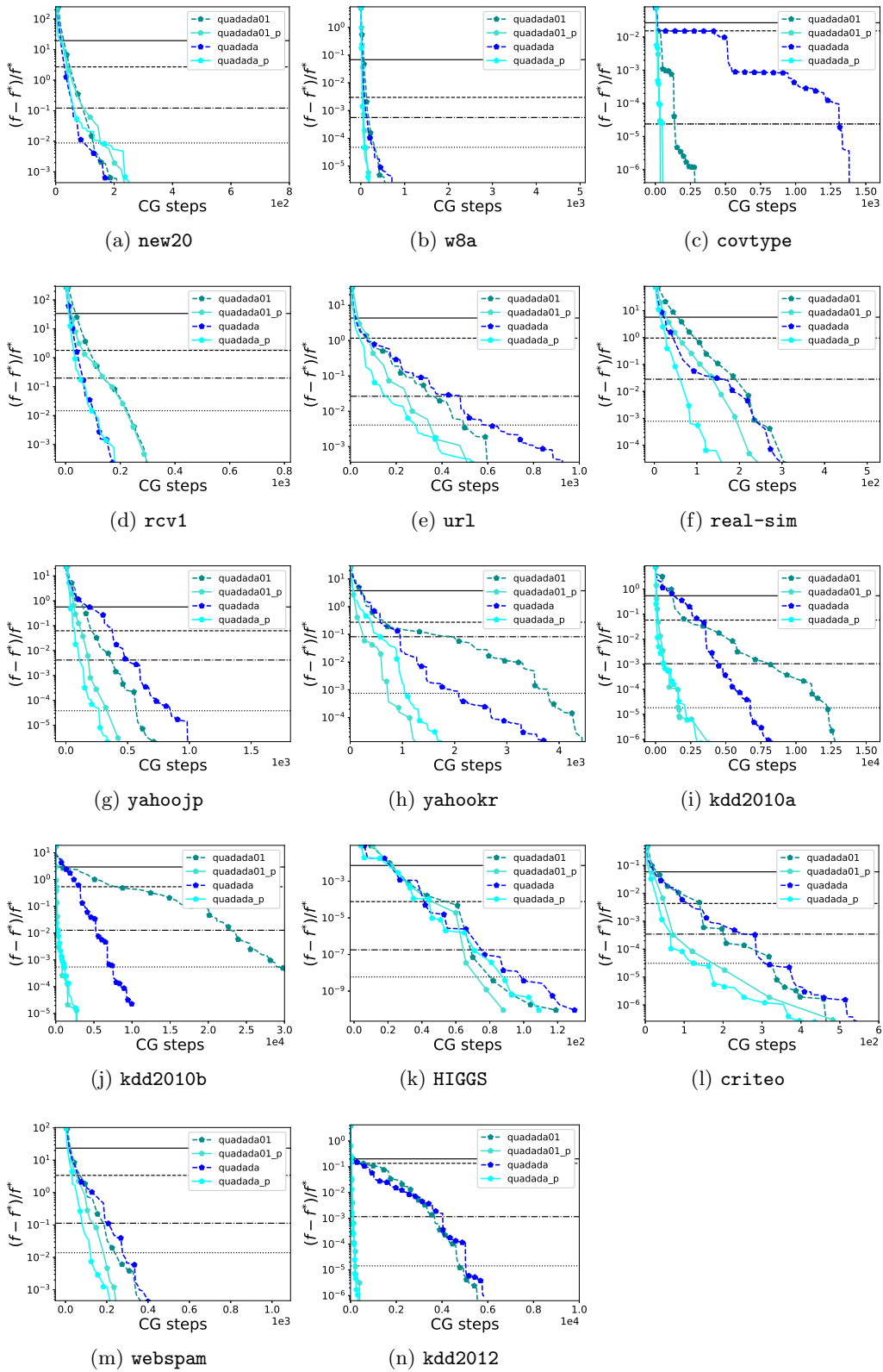


Figure xiv: A comparison of different early thresholds in the quadratic adaptive inner stopping condition. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

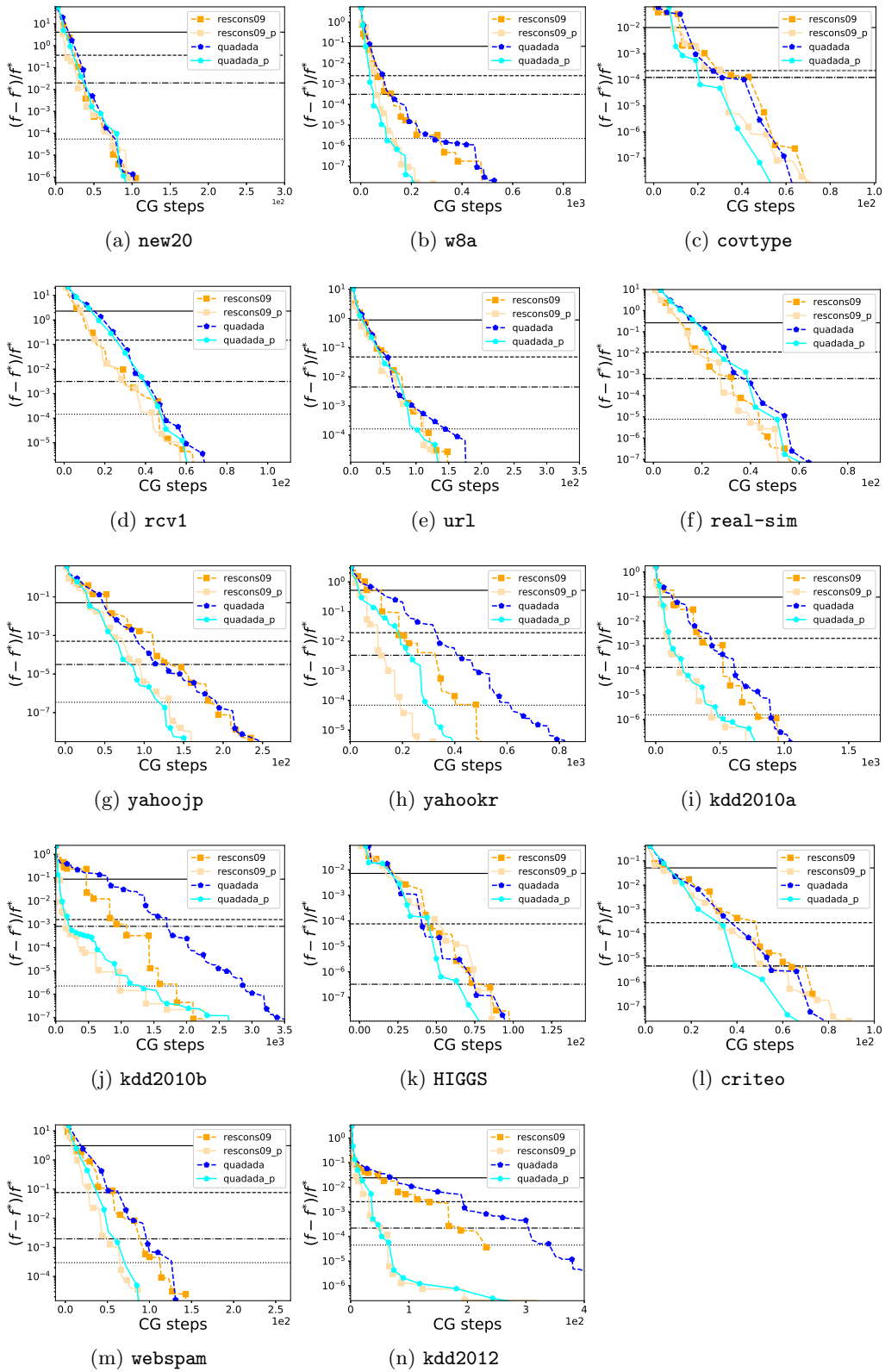


Figure xv: Comparison against a residual constant rule with higher threshold in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

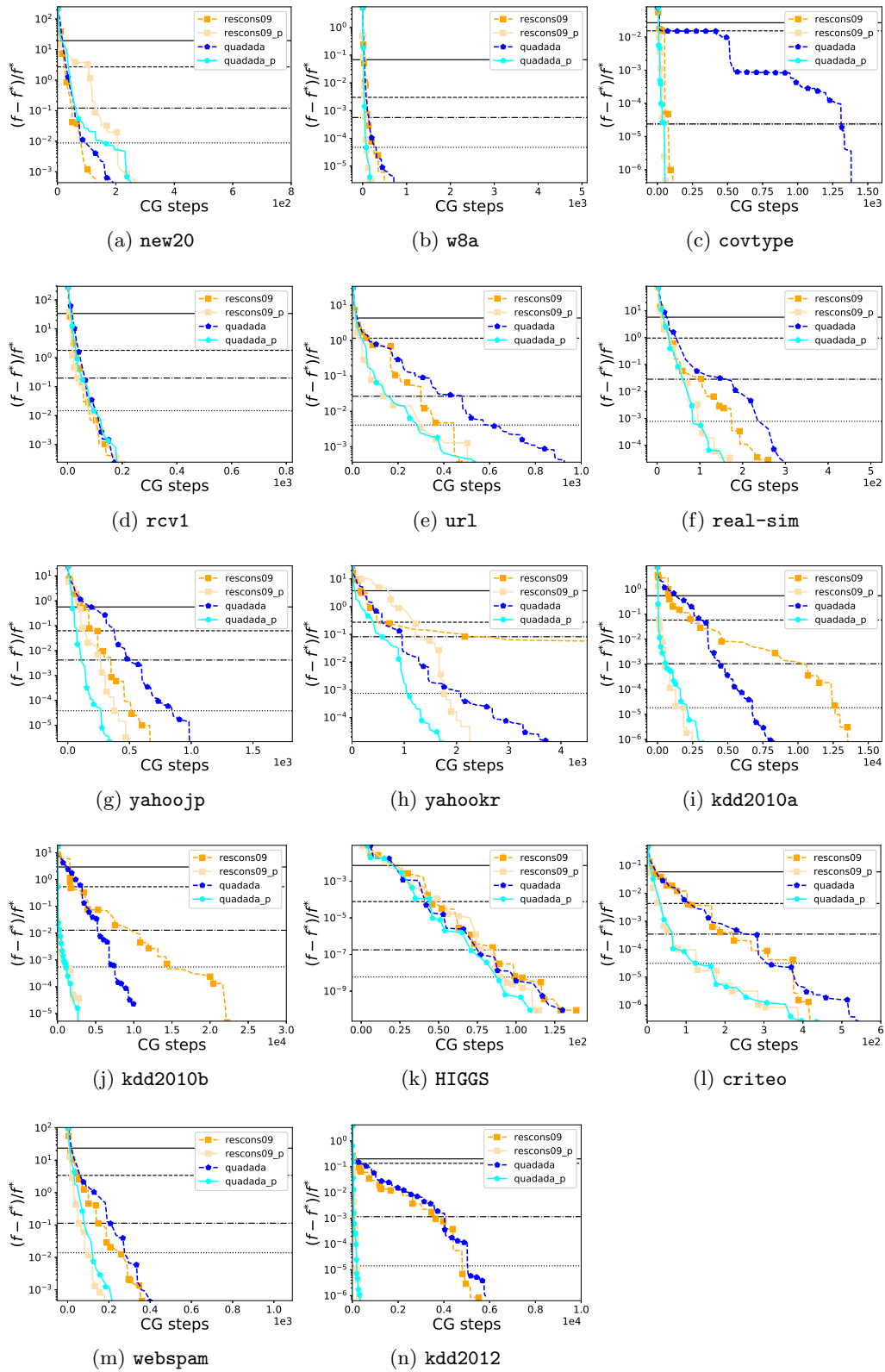


Figure xvi: Comparison against a residual constant rule with higher threshold in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

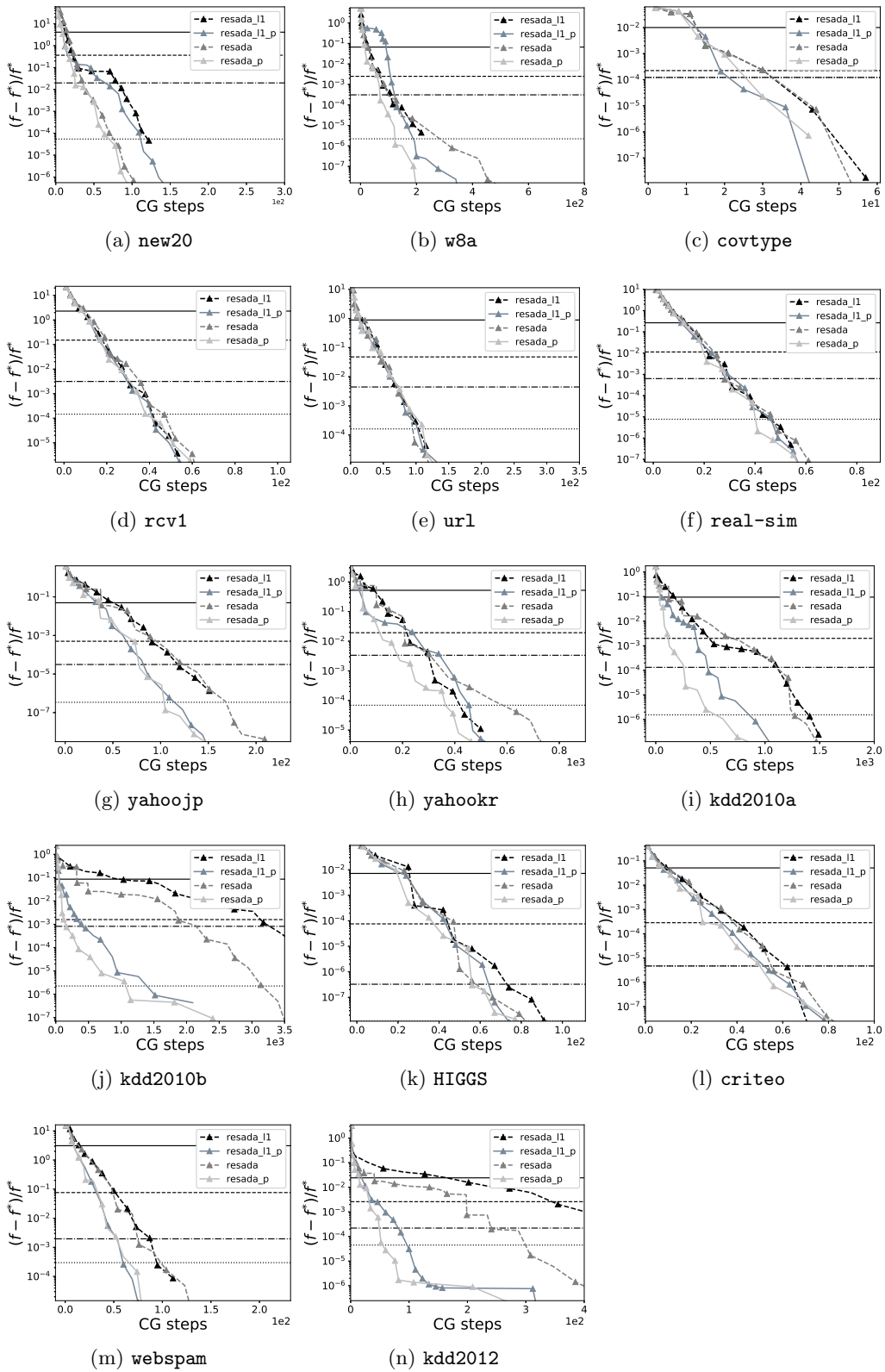


Figure xvii: An investigation on the robustness of residual adaptive rules in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = C_{\text{Best}}$ in each sub-figure.

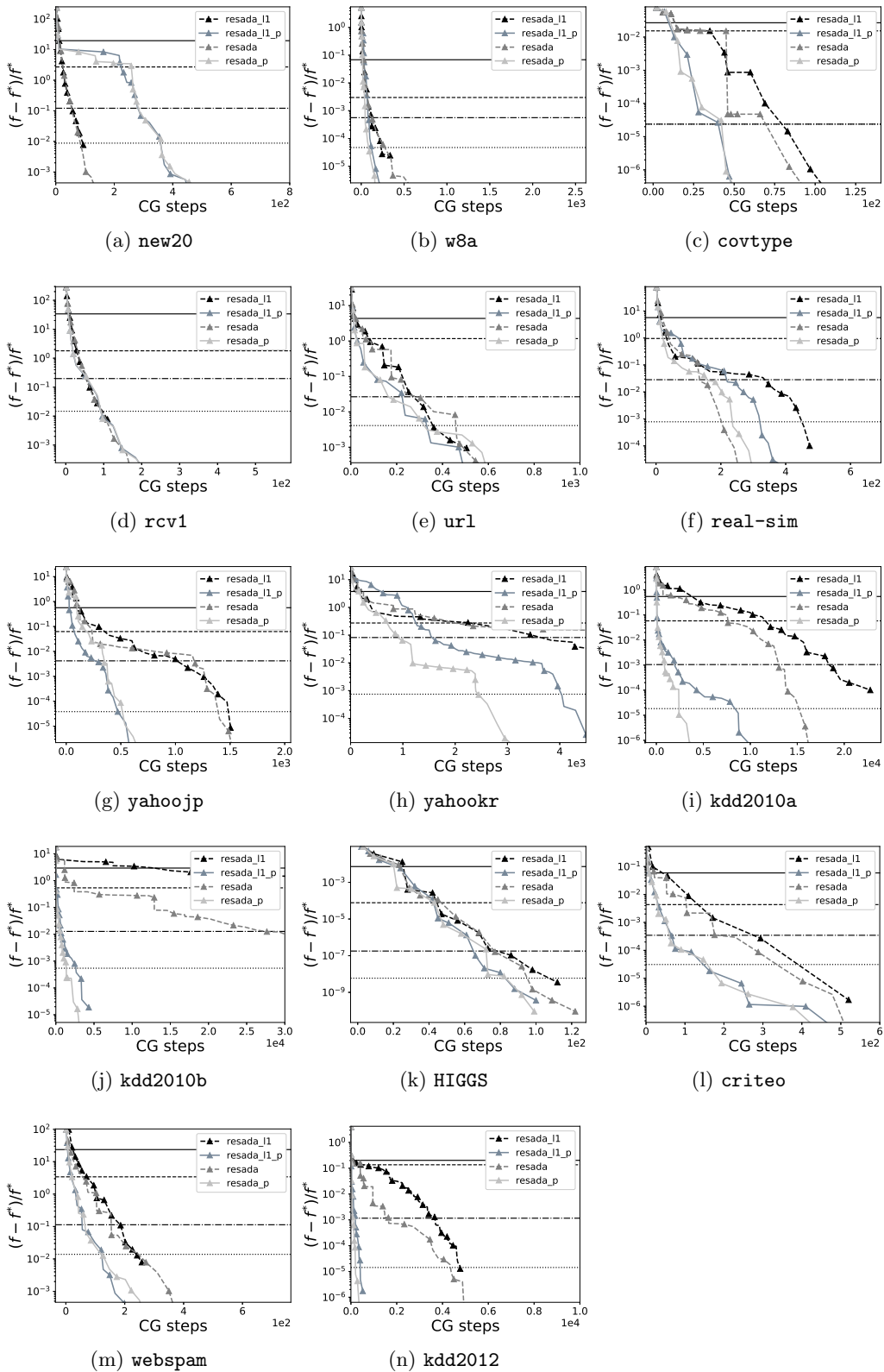


Figure xviii: An investigation on the robustness of residual adaptive rules in the preconditioned case. We show the convergence of a truncated Newton method for logistic regression. The x-axis shows the cumulative number of CG steps. Loss=LR and $C = 100C_{\text{Best}}$ in each sub-figure.

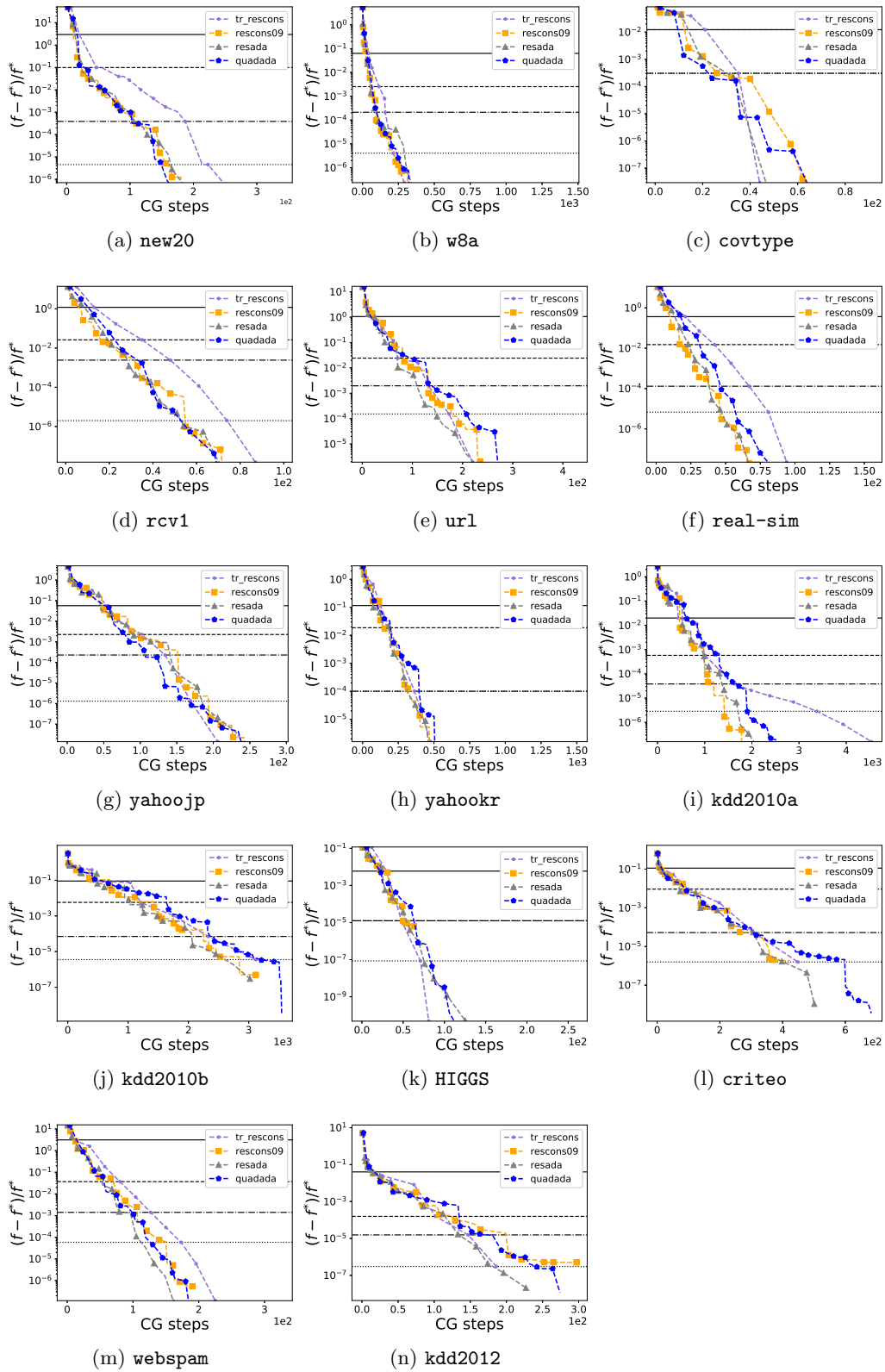


Figure xix: A summary comparison between various forcing sequences. We show the convergence of a truncated Newton method for linear SVM. The x-axis shows the cumulative number of CG steps. Loss=L2 and $C = C_{\text{Best}}$ in each sub-figure.

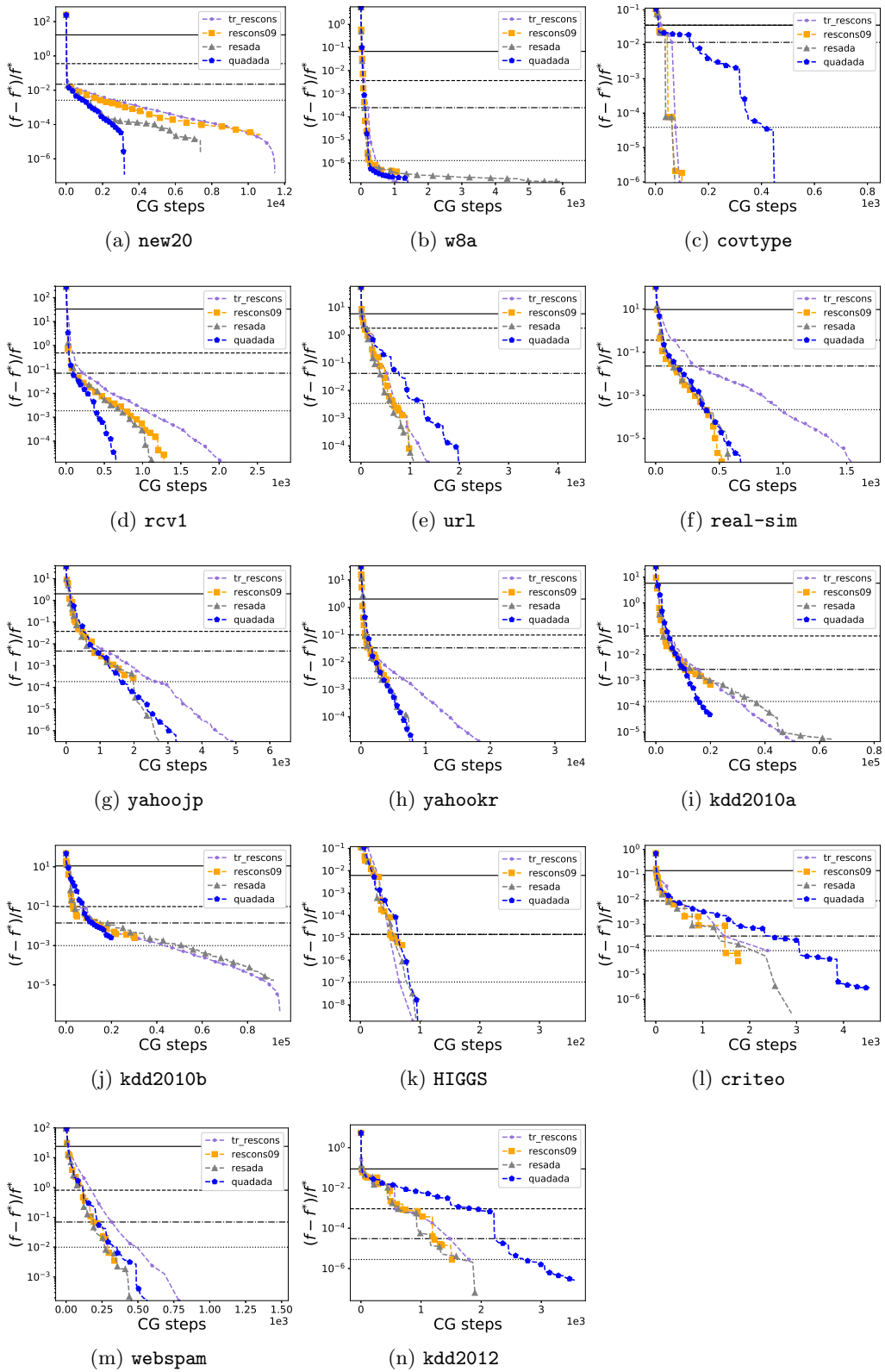


Figure xx: A summary comparison between various forcing sequences. We show the convergence of a truncated Newton method for linear SVM. The x-axis shows the cumulative number of CG steps. Loss=L2 and $C = 100C_{\text{Best}}$ in each sub-figure.

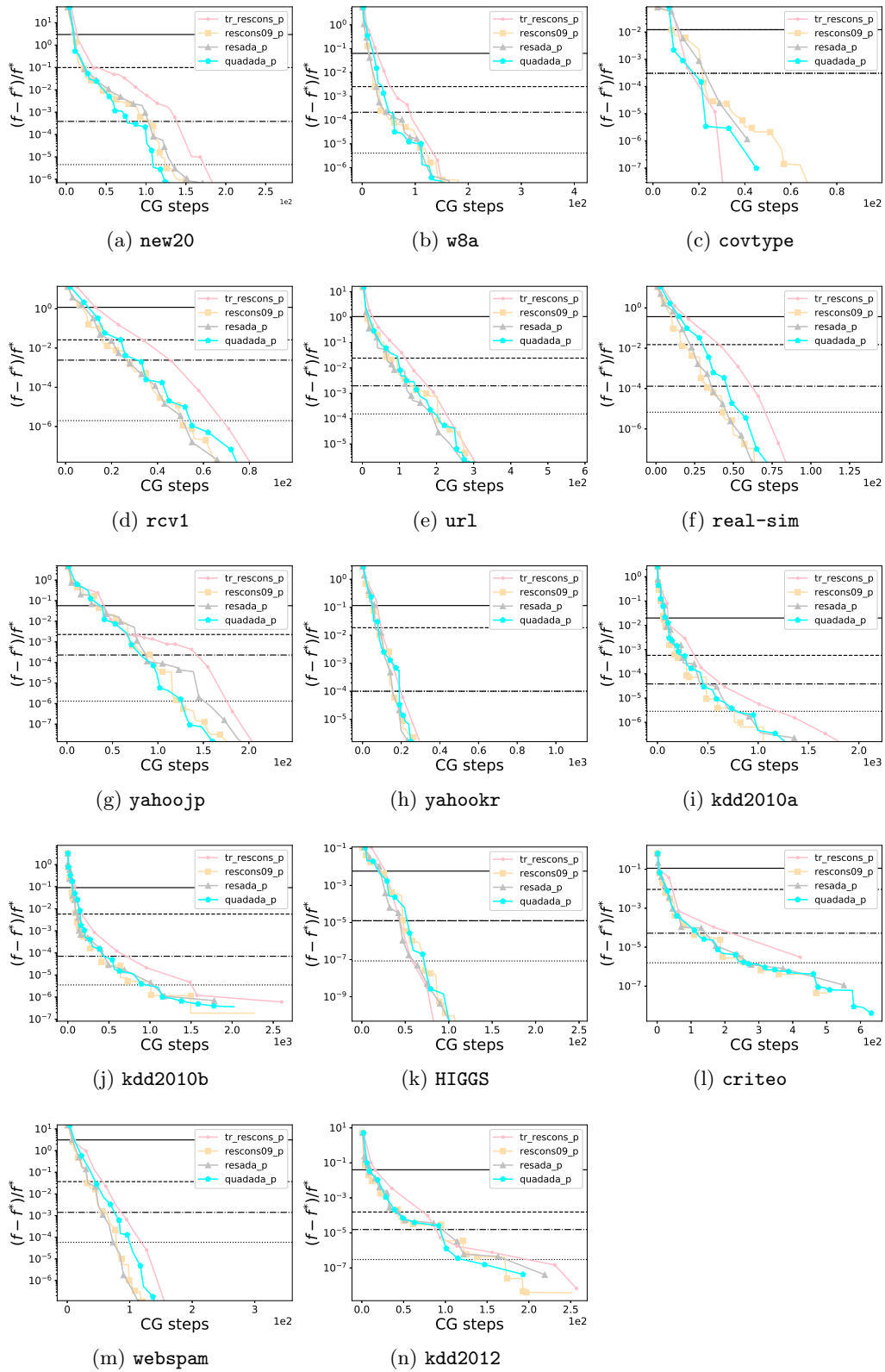


Figure xxi: A summary comparison between various forcing sequences in the preconditioned case. We show the convergence of a truncated Newton method for linear SVM. The x-axis shows the cumulative number of CG steps. Loss=L2 and $C = C_{\text{Best}}$ in each sub-figure.

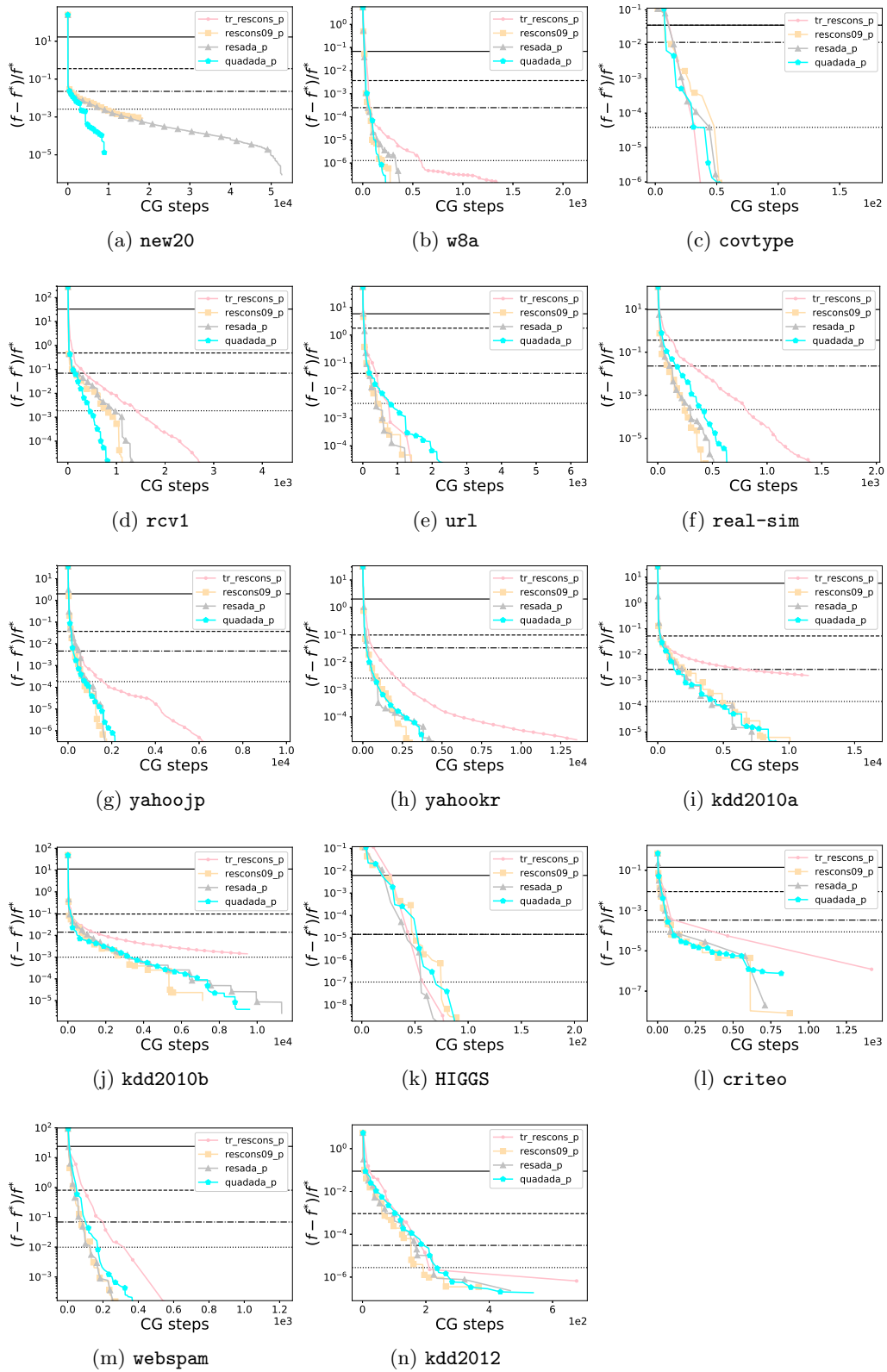


Figure xxii: A summary comparison between various forcing sequences in the preconditioned case. We show the convergence of a truncated Newton method for linear SVM. The x-axis shows the cumulative number of CG steps. Loss=L2 and $C = 100C_{\text{Best}}$ in each sub-figure.