

### 358 A Omitted proofs from Section 3

359 In Section 3, we stated Lemma 9 and proved the first part of it (a moment bound for the determinant).  
 360 Here, we provide the proof of the second part (a moment bound for the adjugate).

361 **Lemma 11 (Lemma 9b restated)** *Let  $\mathbf{A} = \frac{1}{\gamma} \sum_i b_i \mathbf{Z}_i + \mathbf{B}$ , where  $b_i \sim \text{Bernoulli}(\gamma)$  are inde-*  
 362 *pendent, whereas  $\mathbf{Z}_i$  and  $\mathbf{B}$  are  $d \times d$  psd matrices such that  $\|\mathbf{Z}_i\| \leq \epsilon$  for all  $i$  and  $\mathbb{E}[\mathbf{A}] = \mathbf{I}$ . If*  
 363  *$\gamma \geq 8\epsilon d\eta^{-2}(p + \ln d)$  for  $0 < \eta \leq 0.25$  and  $p \geq 2$ , then*

$$\mathbb{E}[\|\text{adj}(\mathbf{A}) - \mathbf{I}\|^p]^{\frac{1}{p}} \leq 9\eta.$$

364 **Proof** Let  $\lambda_{\max}$  and  $\lambda_{\min}$  denote the largest and smallest eigenvalue of  $\text{adj}(\mathbf{A})$ . We have

$$\begin{aligned} \mathbb{E}[\|\text{adj}(\mathbf{A}) - \mathbf{I}\|^p] &= \int_0^\infty px^{p-1} \Pr(\|\text{adj}(\mathbf{A}) - \mathbf{I}\| \geq x) dx \\ &\leq \eta^p + \int_\eta^\infty px^{p-1} (\Pr(\lambda_{\max} \geq 1+x) + \Pr(\lambda_{\min} \leq 1-x)) dx. \end{aligned}$$

365 We will now bound the two probabilities. Let  $\delta_{\max}$  and  $\delta_{\min}$  denote the largest and smallest eigenvalue  
 366 of matrix  $\mathbf{A} - \mathbf{I}$ . Recall the following concentration bounds implied by Lemma 8 (see the first part of  
 367 the proof of Lemma 9):

$$\max \left\{ \Pr(\text{tr}(\mathbf{A} - \mathbf{I}) \geq y), \Pr(\text{tr}(\mathbf{A} - \mathbf{I}) \leq -y) \right\} \leq \begin{cases} e^{-y^2 \frac{2p}{\eta^2}} & \text{for } y \in [0, d]; \\ e^{-y \frac{2dp}{\eta^2}} & \text{for } y \geq d, \end{cases} \quad (6)$$

$$\max \left\{ \Pr(\delta_{\max} \geq z), \Pr(\delta_{\min} \leq -z) \right\} \leq \begin{cases} e^{-z^2 \frac{2p}{\eta^2}} & \text{for } z \in [0, \frac{\eta}{\sqrt{2d}}]; \\ e^{-z^2 \frac{2dp}{\eta^2}} & \text{for } z \in [\frac{\eta}{\sqrt{2d}}, 1]; \\ e^{-z \frac{2dp}{\eta^2}} & \text{for } z \geq 1. \end{cases} \quad (7)$$

368 From the formula  $\text{adj}(\mathbf{A}) = \det(\mathbf{A})\mathbf{A}^{-1}$  it follows that  $\lambda_{\max} \leq \frac{\det(\mathbf{A})}{1+\delta_{\min}} \leq \frac{e^{\text{tr}(\mathbf{A}-\mathbf{I})}}{1+\delta_{\min}}$  so we have

$$\begin{aligned} \Pr(\lambda_{\max} \geq 1+x) &\leq \Pr\left(\frac{e^{\text{tr}(\mathbf{A}-\mathbf{I})}}{1+\delta_{\min}} \geq 1+x\right) \\ &= \Pr\left(\text{tr}(\mathbf{A} - \mathbf{I}) + \ln \frac{1}{1+\delta_{\min}} \geq \ln(1+x)\right) \\ &\leq \Pr\left(\text{tr}(\mathbf{A} - \mathbf{I}) \geq \frac{2}{3} \cdot \ln(1+x)\right) + \Pr\left(\ln \frac{1}{1+\delta_{\min}} \geq \frac{1}{3} \cdot \ln(1+x)\right) \\ &= \Pr\left(\text{tr}(\mathbf{A} - \mathbf{I}) \geq \frac{2}{3} \cdot \ln(1+x)\right) + \Pr\left(\delta_{\min} \leq \frac{1}{(1+x)^{\frac{1}{3}}} - 1\right). \\ &\leq \begin{cases} e^{-\ln^2(1+x) \frac{8p}{9\eta^2}} + e^{-(1-(\frac{1}{1+x})^{\frac{1}{3}})^2 \frac{2p}{\eta^2}} \leq 2e^{-x^2 \frac{p}{20\eta^2}} & \text{for } x \in [0, e-1], \\ e^{-\ln(1+x) \frac{4p}{3\eta^2}} + e^{-\frac{1}{16} \frac{2dp}{\eta^2}} \leq 2e^{-\ln(1+x) \frac{p}{8\eta^2}} & \text{for } x \in [e-1, e^d-1]. \end{cases} \end{aligned}$$

369 For  $x \geq e^d - 1$ , since  $\lambda_{\max} \leq (1+\delta_{\max})^d \leq e^{d\delta_{\max}}$  and  $\ln(1+x) \geq d$ , we have:

$$\Pr(\lambda_{\max} \geq 1+x) \leq \Pr(e^{d\delta_{\max}} \geq 1+x) = \Pr(\delta_{\max} \geq \ln(1+x)/d) \leq e^{-\ln(1+x) \frac{2p}{\eta^2}}.$$

370 Next, we use the fact that for  $\delta = \max\{|\delta_{\max}|, |\delta_{\min}|\}$  we have:

$$\lambda_{\min} \geq \frac{\det(\mathbf{A})}{1+\delta_{\max}} \geq \frac{(1-\delta^2)^d}{(1+\delta_{\max})e^{\text{tr}(\mathbf{I}-\mathbf{A})}} \geq (1-\delta)(1-d\delta^2)(1-\text{tr}(\mathbf{I}-\mathbf{A})),$$

371 so for  $x \in [\eta, 1]$  we have:

$$\begin{aligned} \Pr(\lambda_{\min} \leq 1-x) &\leq \Pr(\delta \geq x/3) + \Pr(\delta^2 \geq x/3d) + \Pr(\text{tr}(\mathbf{I}-\mathbf{A}) \geq x/3) \\ &\leq 2e^{-x^2 \frac{2dp}{9\eta^2}} + 2e^{-\frac{x}{3d} \frac{2dp}{\eta^2}} + e^{-x^2 \frac{2p}{\eta^2}} \leq 5e^{-x^2 \frac{2p}{9\eta^2}}. \end{aligned}$$

Putting everything together we obtain that:

$$\begin{aligned}\mathbb{E}[\|\text{adj}(\mathbf{A}) - \mathbf{I}\|^p] &\leq \eta^p + \int_{\eta}^{e-1} px^{p-1} 7e^{-x^2 \frac{p}{20\eta^2}} dx + \int_{e-1}^{\infty} px^{p-1} 3e^{-\ln(1+x) \frac{p}{8\eta^2}} dx \\ &\leq \eta^p + 7\sqrt{20\pi p} \eta^p + \frac{3p}{\frac{p}{16\eta^2} - 1} \left(\frac{1}{2}\right)^{\frac{p}{16\eta^2} - 1} \\ &\leq \eta^p + 7\sqrt{20\pi p} \eta^p + 6(3\eta)^p \leq (9\eta)^p,\end{aligned}$$

which completes the proof. ■

374

As a consequence of the moment bounds shown in Lemma 9, we establish convergence with high probability for the average of determinants and the adjugates. For the adjugate matrix, we require a matrix variant of the Khintchine/Rosenthal inequalities.

**Lemma 12 ([GCT12])** Suppose that  $p \geq 2$  and  $r = \max\{p, 2 \log d\}$ .<sup>2</sup> Consider a finite sequence  $\{\mathbf{X}_i\}$  of independent, symmetrically random, self-adjoint matrices with dimension  $d \times d$ . Then,

$$\mathbb{E}\left[\left\|\sum_i \mathbf{X}_i\right\|^p\right]^{\frac{1}{p}} \leq \sqrt{er} \left\|\sum_i \mathbb{E}[\mathbf{X}_i^2]\right\|^{\frac{1}{2}} + 2er \mathbb{E}[\max_i \|\mathbf{X}_i\|^p]^{\frac{1}{p}}.$$

**Corollary 13 (Corollary 10 restated)** There is  $C > 0$  s.t. for  $\mathbf{A}$  as in Lemma 9 with all  $\mathbf{Z}_i$  rank-1 and  $\gamma \geq C\epsilon d\eta^{-2} \log^3 \frac{d}{\delta}$ ,

$$(a) \Pr\left(\left|\frac{1}{m} \sum_{t=1}^m \det(\mathbf{A}_t) - 1\right| \geq \frac{\eta}{\sqrt{m}}\right) \leq \delta \quad \text{and} \quad (b) \Pr\left(\left\|\frac{1}{m} \sum_{t=1}^m \text{adj}(\mathbf{A}_t) - \mathbf{I}\right\| \geq \frac{\eta}{\sqrt{m}}\right) \leq \delta,$$

where  $\mathbf{A}_1, \dots, \mathbf{A}_m$  are independent copies of  $\mathbf{A}$ .

**Proof** Applying Lemma 9 to the matrix  $\mathbf{A}$ , for appropriate  $C$  and any fixed  $p \geq 2$ , if  $\gamma \geq C\epsilon d\sigma^{-2}(p + \ln d)$ , then for any  $s \in [2, p]$  we have  $\mathbb{E}[\|\text{adj}(\mathbf{A}_t) - \mathbf{I}\|^s] \leq \sigma^s$ . With the additional assumption that  $\mathbf{Z}_i$ 's are rank-1, Theorem 7 implies that  $\mathbb{E}[\text{adj}(\mathbf{A}_t)] = \mathbf{I}$ , so by a standard symmetrization argument, where  $r_t$  denote independent Rademacher random variables,

$$\mathbb{E}\left[\left\|\frac{1}{m} \sum_{t=1}^m \text{adj}(\mathbf{A}_t) - \mathbf{I}\right\|^p\right]^{\frac{1}{p}} \leq 2 \cdot \mathbb{E}\left[\left\|\sum_t \frac{r_t}{m} (\text{adj}(\mathbf{A}_t) - \mathbf{I})\right\|^p\right]^{\frac{1}{p}}.$$

Applying Lemma 12 to the matrices  $\mathbf{X}_t = \frac{1}{m} \mathbf{Y}_t$ , where  $\mathbf{Y}_t = r_t (\text{adj}(\mathbf{A}_t) - \mathbf{I})$ , we obtain that:

$$\begin{aligned}\mathbb{E}\left[\left\|\frac{1}{m} \sum_t \mathbf{Y}_t\right\|^p\right]^{\frac{1}{p}} &\leq \sqrt{er} \left\|m \cdot \frac{1}{m^2} \mathbb{E}[\mathbf{Y}_t^2]\right\|^{\frac{1}{2}} + \frac{2er}{m} \mathbb{E}\left[\sum_{i=1}^m \|\mathbf{Y}_i\|^p\right]^{\frac{1}{p}} \\ &\leq \sqrt{\frac{er}{m}} \cdot \mathbb{E}[\|\mathbf{Y}\|^2]^{\frac{1}{2}} + \frac{2er}{m} (m \cdot \mathbb{E}[\|\mathbf{Y}\|^p])^{\frac{1}{p}} \\ &\leq \left(\sqrt{\frac{er}{m}} + \frac{2er}{m^{1-\frac{1}{p}}}\right) \cdot \sigma \leq C' \cdot \frac{p\sigma}{\sqrt{m}},\end{aligned}$$

for  $p \geq 2 \log d$  and  $C'$  chosen appropriately. Now Markov's inequality yields:

$$\Pr\left(\left\|\frac{1}{m} \sum_t \text{adj}(\mathbf{A}_t) - \mathbf{I}\right\| \geq \alpha\right) \leq \alpha^{-p} \cdot \mathbb{E}\left[\left\|\frac{1}{m} \sum_t \text{adj}(\mathbf{A}_t) - \mathbf{I}\right\|^p\right] \leq \left(\frac{2C'p\sigma}{\alpha\sqrt{m}}\right)^p.$$

Setting  $\alpha = \frac{\eta}{\sqrt{m}}$ ,  $\sigma = \frac{\eta}{4C'p}$  and  $p = 2 \lceil \max\{\log d, \log \frac{1}{\delta}\} \rceil$ , the above bound becomes  $(\frac{1}{2})^p \leq \delta$  for  $k \geq C'' \mu d^2 \eta^{-2} (\log^3 \frac{1}{\delta} + \log^3 d)$ . Showing the analogous result for the average of determinants of matrices  $\mathbf{A}_t$  instead of the adjugates follows identically, except that Lemma 12 can be replaced with the standard scalar Rosenthal's inequality. ■

393

<sup>2</sup>In [GCT12] it is assumed that  $d \geq 3$ , however this assumption is not used anywhere in the proof.

## 394 B Proof of Newton convergence

395 Here, we provide a proof of Corollary 6, which describes the convergence guarantees for the  
 396 approximate Newton step obtained via determinantal averaging. It suffices to show the following  
 397 lemma.

398 **Lemma 14** *Let loss  $\mathcal{L}$  be defined as in (1) and assume its Hessian is  $L$ -Lipschitz (Assumption 5). If*

$$\|\hat{\mathbf{p}} - \mathbf{p}^*\|_{\nabla^2 \mathcal{L}(\mathbf{w})} \leq \alpha \|\mathbf{p}^*\|_{\nabla^2 \mathcal{L}(\mathbf{w})}, \quad \text{where } \mathbf{p}^* = \nabla^{-2} \mathcal{L}(\mathbf{w}) \nabla \mathcal{L}(\mathbf{w}),$$

399 *then the approximate Newton step  $\tilde{\mathbf{w}} = \mathbf{w} - \hat{\mathbf{p}}$  satisfies:*

$$\|\tilde{\mathbf{w}} - \mathbf{w}^*\| \leq \max \left\{ \alpha \sqrt{\kappa} \|\mathbf{w} - \mathbf{w}^*\|, \frac{2L}{\sigma_{\min}} \|\mathbf{w} - \mathbf{w}^*\|^2 \right\}, \quad \text{where } \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w}),$$

400 *where  $\kappa$  and  $\sigma_{\min}$  are the conditioning number and smallest eigenvalue of  $\nabla^2 \mathcal{L}(\mathbf{w})$ , respectively.*

401 **Proof** The lemma essentially follows via the standard analysis of the Newton's method. For the  
 402 sake of completeness we will outline the proof following [WRXM17]. Denoting  $\mathbf{H} = \nabla^2 \mathcal{L}(\mathbf{w})$  and  
 403  $\mathbf{g} = \nabla \mathcal{L}(\mathbf{w})$ , we define the auxiliary function

$$\phi(\mathbf{p}) \stackrel{\text{def}}{=} \mathbf{p}^\top \mathbf{H} \mathbf{p} - 2\mathbf{p}^\top \mathbf{g}$$

404 By definition of  $\phi(\mathbf{p})$  we have  $\phi(\mathbf{p}^*) = \phi(\mathbf{H}^{-1} \mathbf{g}) = -\|\mathbf{p}^*\|_{\mathbf{H}}^2$ . It follows that

$$\begin{aligned} \phi(\hat{\mathbf{p}}) - \phi(\mathbf{p}^*) &= \|\mathbf{H}^{\frac{1}{2}} \hat{\mathbf{p}}\|^2 - 2\mathbf{g}^\top \mathbf{H}^{-1} \mathbf{H} \hat{\mathbf{p}} + \|\mathbf{H}^{\frac{1}{2}} \mathbf{p}^*\|^2 \\ &= \|\mathbf{H}^{\frac{1}{2}} (\hat{\mathbf{p}} - \mathbf{p}^*)\|^2 = \|\hat{\mathbf{p}} - \mathbf{p}^*\|_{\mathbf{H}}^2 \leq \alpha^2 \|\mathbf{p}^*\|_{\mathbf{H}}^2 = -\alpha^2 \phi(\mathbf{p}^*). \end{aligned}$$

405 We invoke the classical result in local convergence analysis of Newton's method [NW06], using the  
 406 statement of Lemma 9 in [WRXM17].

407 **Lemma 15 ([WRXM17])** *Assume Hessian is  $L$ -Lipschitz and that  $\hat{\mathbf{p}}$  satisfies  $\phi(\hat{\mathbf{p}}) \leq (1 -$   
 408  $\alpha^2) \min_{\mathbf{p}} \phi(\mathbf{p})$ . Then  $\tilde{\mathbf{w}} = \mathbf{w} - \hat{\mathbf{p}}$  satisfies*

$$\|\tilde{\mathbf{w}} - \mathbf{w}^*\|_{\mathbf{H}}^2 \leq L \|\mathbf{w} - \mathbf{w}^*\|^2 \|\tilde{\mathbf{w}} - \mathbf{w}^*\| + \frac{\alpha^2}{1 - \alpha^2} \|\mathbf{w} - \mathbf{w}^*\|_{\mathbf{H}}^2.$$

409 Lemma 15 immediately implies that one of the following two inequalities hold:

$$\begin{aligned} \|\tilde{\mathbf{w}} - \mathbf{w}^*\| &\leq \frac{2L}{\sigma_{\min}(\mathbf{H})} \cdot \|\mathbf{w} - \mathbf{w}^*\|^2, \\ \|\tilde{\mathbf{w}} - \mathbf{w}^*\| &\leq \frac{\alpha}{\sqrt{1 - \alpha^2}} \sqrt{\frac{2\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})}} \cdot \|\mathbf{w} - \mathbf{w}^*\|, \end{aligned}$$

410 which proves Lemma 14. ■

411 Note that Corollary 6 follows immediately by combining Corollary 4 with Lemma 14.

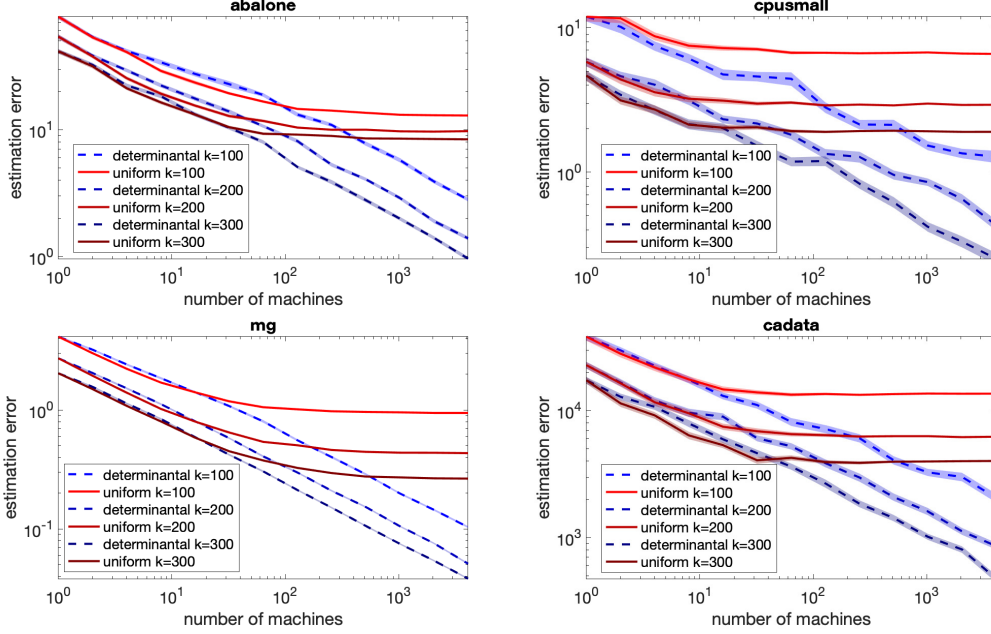


Figure 2: Comparison of the estimation error between *determinantal* and *uniform* averaging on four libsvm datasets.

## C Experiments

In this section, we experimentally evaluate the estimation error of determinantal averaging for the Newton’s method (following the setup of Section 1.1), and we compare it against uniform averaging [WRXM17]. We use square loss  $\ell_i(\mathbf{w}^\top \mathbf{x}_i) = (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$ , where  $y_i$  are the real-valued labels for a regression problem, and we run the experiments on several benchmark regression datasets from the libsvm repository [CL11]. In this setting, the local Newton estimate computed from the starting vector  $\mathbf{w} = \mathbf{0}$  is given by:

$$\hat{\mathbf{p}} = \left( \frac{1}{k} \sum_{i=1}^n b_i \mathbf{x}_i \mathbf{x}_i^\top + \lambda \mathbf{I} \right)^{-1} \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i, \quad \text{where } b_i \sim \text{Bernoulli}(k/n).$$

In all of our experiments we set the regularization parameter to  $\lambda = \frac{1}{n}$ . Let  $\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_m \stackrel{\text{i.i.d.}}{\sim} \hat{\mathbf{p}}$  be  $m$  distributed local estimates and denote  $\hat{\mathbf{H}}_t$  as the  $t$ th local Hessian estimate. The two averaging strategies we compare are:

$$\text{determinantal: } \hat{\mathbf{p}}_{\text{det}} = \frac{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t) \hat{\mathbf{p}}_t}{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t)}, \quad \text{uniform: } \hat{\mathbf{p}}_{\text{uni}} = \frac{1}{m} \sum_{t=1}^m \hat{\mathbf{p}}_t.$$

Figure 2 plots the estimation errors  $\|\hat{\mathbf{p}}_{\text{det}} - \mathbf{p}^*\|$  and  $\|\hat{\mathbf{p}}_{\text{uni}} - \mathbf{p}^*\|$ , where  $\mathbf{p}^*$  is the exact Newton step starting from  $\mathbf{w} = \mathbf{0}$ , for datasets ABALONE, CPUSMALL, MG<sup>3</sup> and CADATA [CL11] (for convenience, the plot from Figure 1 in Section 1.1 is repeated here). The reported results are averaged over 100 trials, with shading representing standard error. We consistently observe that for a small number of machines  $m$  both methods effectively reduce the estimation error, however after a certain point uniform averaging converges to a biased estimate and the estimation error flattens out. On the other hand, determinantal averaging continues to converge to the optimum as the number of machines keeps growing. We remark that for some datasets determinantal averaging exhibits larger variance than uniform averaging, especially when local sample size is small. Reducing that variance, for example through some form of additional regularization, is a new direction for future work.

<sup>3</sup>We expanded features to all degree 2 monomials, and removed redundant ones.