

## Distributed Newton's method

**Task:** Minimization of a convex loss:

$$\mathcal{L}(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^\top \mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{for } \mathbf{w} \in \mathbb{R}^d.$$

**Goal:** Find a good descent direction:  $\tilde{\mathbf{w}} = \mathbf{w} - \mathbf{p}$

Newton's method: use both *Hessian* and *gradient* information,

$$\mathbf{p} = \mathbf{H}^{-1} \mathbf{g}, \quad \text{where } \mathbf{H} = \nabla^2 \mathcal{L}(\mathbf{w}), \quad \mathbf{g} = \nabla \mathcal{L}(\mathbf{w}).$$

*Distributed Newton:* Avoid constructing the full Hessian by replacing it with local approximations computed on separate machines:

$$\hat{\mathbf{p}}_t = \left[ \underbrace{\nabla^2 \hat{\mathcal{L}}_t(\mathbf{w})}_{\text{local Hessian } \hat{\mathbf{H}}_t} \right]^{-1} \underbrace{\nabla \mathcal{L}(\mathbf{w})}_{\text{global gradient } \mathbf{g}} \quad \text{for } t = 1, \dots, m,$$

where  $\hat{\mathcal{L}}$  is based on a random sample of data,

$$\hat{\mathcal{L}}(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{k} \sum_{i=1}^n b_i \ell_i(\mathbf{w}^\top \mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2, \quad \text{where } b_i \sim \text{Bernoulli}(k/n).$$

**Question:** How to combine local Newton estimates  $\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_m$ ?

## Problem: Inversion bias

Standard averaging leads to biased estimates:

$$\lim_{m \rightarrow \infty} \frac{1}{m} \sum_{t=1}^m \hat{\mathbf{p}}_t \neq \mathbf{p} \quad (m \text{ is the number of machines})$$

For large  $m$ , adding more machines will not improve the accuracy

The reason for this is a general phenomenon, which we call *inversion bias*:

$$\mathbb{E}[\hat{\mathbf{H}}^{-1}] \neq \mathbf{H}^{-1}, \quad \text{even though } \mathbb{E}[\hat{\mathbf{H}}] = \mathbf{H}.$$

## Other examples of inversion bias

Consider a data covariance matrix:  $\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ .

In *uncertainty quantification* we wish to estimate:

- the trace of  $\Sigma^{-1}$ ,
- a subset of entries of  $\Sigma^{-1}$ .

Again, we encounter inversion bias when averaging estimates.

## Determinantal averaging

**Goal:** Estimate a linear function of inverse Hessian,  $F(\mathbf{H}^{-1})$

**Given:**  $m$  independent local estimates  $F(\hat{\mathbf{H}}_1), \dots, F(\hat{\mathbf{H}}_m)$

*Newton estimates:*  $F(\hat{\mathbf{H}}_t^{-1}) = \hat{\mathbf{H}}_t^{-1} \mathbf{g} = \hat{\mathbf{p}}_t$ , where  $\mathbf{g}$  is the gradient.

**Strategy:** Weighted average of the estimates,

$$\hat{F}_m = \frac{\sum_{t=1}^m a_t F(\hat{\mathbf{H}}_t^{-1})}{\sum_{t=1}^m a_t}$$

Uniform averaging ( $a_t = \frac{1}{n}$ ) suffers from *inversion bias*:  $\hat{F}_m \not\rightarrow F(\mathbf{H}^{-1})$

**Determinantal averaging:** use carefully-chosen non-uniform weights

$$a_t = \det(\hat{\mathbf{H}}_t) \Rightarrow \text{no inversion bias!}$$

## Main result: Distributed Newton without inversion bias

### Theorem

If expected local sample size satisfies  $k \geq C\eta^{-2}\mu d^2 \log^3 \frac{d}{\delta}$  then

$$\left\| \frac{\sum_{t=1}^m a_t \hat{\mathbf{p}}_t}{\sum_{t=1}^m a_t} - \mathbf{p} \right\|_{\mathbf{H}} \leq \frac{\eta}{\sqrt{m}} \cdot \|\mathbf{p}\|_{\mathbf{H}} \quad \text{with probability } \geq 1 - \delta,$$

where  $\mu = \frac{1}{d} \max_i \ell_i''(\mathbf{w}^\top \mathbf{x}_i) \|\mathbf{x}_i\|_{[\nabla^2 \mathcal{L}(\mathbf{w})]^{-1}}^2$  and  $a_t = \det(\hat{\mathbf{H}}_t)$ .

## Asymptotically consistent inverse estimator

Determinantal averaging is asymptotically consistent:

$$\lim_{m \rightarrow \infty} \frac{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t) F(\hat{\mathbf{H}}_t^{-1})}{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t)} = F(\mathbf{H}^{-1}).$$

Adding more estimates always improves the accuracy

Key underlying expectation formula:

$$\frac{\mathbb{E}[\det(\hat{\mathbf{H}}) \hat{\mathbf{H}}^{-1}]}{\mathbb{E}[\det(\hat{\mathbf{H}})]} = \mathbf{H}^{-1}.$$

### Corollary

A convergence result for Distributed Newton:

$$\|\tilde{\mathbf{w}} - \mathbf{w}^*\| \leq \max \left\{ \frac{\eta}{\sqrt{m}} \sqrt{\kappa} \|\mathbf{w} - \mathbf{w}^*\|, \frac{2L}{\lambda_{\min}} \|\mathbf{w} - \mathbf{w}^*\|^2 \right\}$$

for  $\tilde{\mathbf{w}} = \mathbf{w} - \frac{\sum_{t=1}^m a_t \hat{\mathbf{p}}_t}{\sum_{t=1}^m a_t}$  and  $\mathbf{w}^* = \underset{\mathbf{w}}{\text{argmin}} \mathcal{L}(\mathbf{w})$ .

$L, \kappa, \lambda_{\min}$  - Lipschitz constant, condition number and smallest eigenvalue of  $\mathbf{H}$ .

## Proof techniques

### Key lemma

If  $\hat{\mathbf{H}} = \sum_i s_i \mathbf{Z}_i$ , where  $s_i$  are independent random variables and  $\mathbf{Z}_i$  are fixed square rank-1 matrices, then

$$(a) \quad \underbrace{\mathbb{E}[\det(\hat{\mathbf{H}})] = \det(\mathbb{E}[\hat{\mathbf{H}}])}_{\text{determinant commutes with expectation}} \quad \text{and} \quad (b) \quad \underbrace{\mathbb{E}[\text{adj}(\hat{\mathbf{H}})] = \text{adj}(\mathbb{E}[\hat{\mathbf{H}}])}_{\text{adjugate commutes with expectation}}.$$

*Adjugate matrix:*  $\text{adj}(\mathbf{A}) = \det(\mathbf{A}) \mathbf{A}^{-1}$  for any invertible  $\mathbf{A}$

Main result relies on showing an improved *matrix concentration inequality*:

$$\left(1 - \frac{\eta}{\sqrt{m}}\right) \cdot \mathbf{H}^{-1} \preceq \frac{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t) \hat{\mathbf{H}}_t^{-1}}{\sum_{t=1}^m \det(\hat{\mathbf{H}}_t)} \preceq \left(1 + \frac{\eta}{\sqrt{m}}\right) \cdot \mathbf{H}^{-1}$$

## Experiment

Newton step estimation error versus number of machines  $m$

