

DPPs and GPs

August 6, 2019

1 8/5/2019

1.1 Distributed ridge via reweighting

Dobriban and Sheng (2019) give a distributed model where the ridge estimator is first taken on subsets of data (e.g. in a distributed setting)

$$\hat{\beta}_S := \mathbf{X}_S^\dagger \mathbf{y}_S$$

and the local estimates are averaged

$$\hat{\beta} := \sum_{S_i} w_i \hat{\beta}_{S_i}$$

We propose an estimator consisting of averaged local estimates. Our method is unique in that we construct a globally regularized estimator from unregularized local estimators. This is achieved through exploiting the implicit regularization present in DPPs via importance reweighting.

Proposition 1. Let $S \subset [n]$ and define the local estimate $\hat{\beta}_S := \mathbf{X}_S^\dagger \mathbf{y}_S$.

For $|S| = k \geq d$, if we use weights

$$w_S = \frac{P(S)}{Q(S)}$$

where

$$P(S) = \frac{\det \lambda^{-1} \mathbf{X}_S^\top \mathbf{X}_S}{\sum_{|S|=k} \det \mathbf{X}_S^\top \mathbf{X}_S} = \frac{\det \lambda^{-1} \mathbf{X}_S^\top \mathbf{X}_S}{\binom{n-d}{k-d} \det \mathbf{X}^\top \mathbf{X}}$$

and $Q(S)$ is the probability that the subset $S \subset [n]$ is present on any given machine (assumed IID), then

$$\mathbb{E} \hat{\beta} = \mathbb{E} \frac{1}{|I|} \sum_i w_{S_i} \hat{\beta}_{S_i} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \quad (1)$$

where I is the index set over all machines.

Remark 1. This shows that the averaged estimator $\hat{\beta}$ provides an unbiased estimator to the full ridge solution. An interesting property is that we are averaging ordinary least squares estimators; the regularization is implicit and obtained by the weights w_S .

Remark 2. By weighting estimates obtained on “larger” subsets \mathbf{X}_S , we give more importance to estimators using “larger” amounts of data.

Interpret
as a form
of variance
reduction.

Remark 3. The proposition implies a trivial map-reduce algorithm which (assuming the data is already partitioned) can produce $\hat{\beta}$ in $O(k^3 + |I|D)$ time.

Proof.

$$\begin{aligned}\mathbb{E}\hat{\beta} &= \frac{1}{|I|} \sum_i \mathbb{E} w_{S_i} \hat{\beta}_{S_i} \\ &= \sum_{|S|=k} Q(S) \frac{P(S)}{Q(S)} \hat{\beta}_S \\ &= \sum_{|S|=k} P(S) \hat{\beta}_S\end{aligned}$$

Result follows by expectation lemma for regularized volume sampling, ■

1.2 Follow up from Michal's notes

Proof. (Proof of Lemma 1 from determinantal random projections) By the definition of Fredholm determinant (analogously, by L-ensemble normalizing constant result of Machi)

$$\det(\mathbf{I} + \Sigma_\mu) = \sum_{k=0}^d \sum_{|S|=k} \det(\Sigma_\mu)_S$$

where S ranges over all increasing subsequences of $[n]$ of length k . Hence, it suffices to show $\sum_{|S|=k} \det(\Sigma_\mu)_S = \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim \mu^k} [\det \mathbf{X} \mathbf{Y}^\top]$. Unpacking the summand, we see that

$$\sum_{|S|=k} \det(\Sigma_\mu)_S = \sum_{|S|=k} \det \mathbb{E}[\mathbf{P}_S \mathbf{x} \mathbf{y}^\top \mathbf{P}_S^\top]$$

We would like to simplify the sum over S using Cauchy-Binet. However, this is only possible when the outer dimension is fixed and S is on the inside, i.e.

Here's where I get stuck

$$\sum_{|S|=k} \det \mathbf{X} \mathbf{P}_S \mathbf{P}_S^\top \mathbf{Y} = \binom{d-k}{k-k} \det \mathbf{X} \mathbf{Y}^\top$$

(when we take $\mathbf{X} \in \mathbb{R}^{k \times d}$). ■

Theorem 1 (Counterexample to Lemma 2 in determinantal random projections). *Let μ be the uniform measure over a discrete set $[n]$ indexing the rows of $\mathbf{X} \in \mathbb{R}^{n \times d}$. Then if $S \sim \text{DPP}(L = \mathbf{X} \mathbf{X}^\top)$, we have*

$$\mathbb{E}[\mathbf{f}_S^\top (\mathbf{X}_S \mathbf{X}_S^\top)^{-1} \mathbf{g}_S] = \mathbf{f}^\top (\mathbf{I} + \mathbf{X} \mathbf{X}^\top) \mathbf{g} \quad (2)$$

Proof. First notice that this L-ensemble is the same as the setting in Lemma 2, since for ordered sequences $S \subset [n]$ of size k we have

$$P(S) \propto \frac{1}{k!} \det \mathbf{X}_S \mathbf{X}_S^\top \propto \frac{1}{k!} \mu(S) \det \mathbf{X}_S \mathbf{X}_S^\top = \frac{1}{k!} \mathbb{E}_T \left[\mathbb{1}_{T=S} \det \mathbf{X}_T \mathbf{X}_T^\top \right] = \frac{1}{k!} \mathbb{E}_{\mathbf{Z} \sim \mu^k} \left[\mathbb{1}_{\{z_i\}_{i=1}^k = S} \det \mathbf{Z} \mathbf{Z}^\top \right] \quad (3)$$

Next, note that the resulting conclusion differs:

$$\det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \mathbb{E} \mathbf{f}_S^\top (\mathbf{X}_S \mathbf{X}_S)^{-1} \mathbf{g}_S = \sum_{\substack{S \subset [n] \\ 0 \leq |S| \leq d}} \det(\mathbf{X}_S \mathbf{X}_S^\top) \mathbf{f}_S^\top (\mathbf{X}_S \mathbf{X}_S)^{-1} \mathbf{g}_S \quad (4)$$

$$= \sum_{\substack{S \subset [n] \\ 0 \leq |S| \leq d}} \mathbf{f}_S^\top \text{adj}(\mathbf{X}_S \mathbf{X}_S)^{-1} \mathbf{g}_S \quad (5)$$

$$= \sum_{\substack{S \subset [n] \\ 0 \leq |S| \leq d}} \det(\mathbf{X}_S \mathbf{X}_S^\top + \mathbf{f}_S \mathbf{g}_S^\top) - \det \mathbf{X}_S \mathbf{X}_S^\top \quad (6)$$

$$= \left(\sum_{\substack{S \subset [n] \\ 0 \leq |S| \leq d}} \det([\mathbf{X}, \mathbf{f}]_S [\mathbf{X}, \mathbf{g}]_S^\top) \right) - \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \quad (7)$$

$$= \det([\mathbf{X}, \mathbf{f}] [\mathbf{X}, \mathbf{g}]^\top) - \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \quad (8)$$

$$= \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top + \mathbf{f}\mathbf{g}^\top) - \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \quad (9)$$

$$= \mathbf{f}^\top \text{adj}(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \mathbf{g} + \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) - \det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \quad (10)$$

$$= \mathbf{f}^\top \text{adj}(\mathbf{I} + \mathbf{X}\mathbf{X}^\top) \mathbf{g} \quad (11)$$

Divide both sides by $\det(\mathbf{I} + \mathbf{X}\mathbf{X}^\top)$ to get the desired result. \blacksquare

Proposition 2 (Sylvester's for operators). *For trace class operators $A : H_1 \rightarrow H_2$ and $B : H_2 \rightarrow H_1$ on Hilbert spaces*

$$\det(\mathbf{I} + \mathbf{A}\mathbf{B}) = \det(\mathbf{I} + \mathbf{B}\mathbf{A}) \quad (12)$$

Proof.

$$\det(\mathbf{I} + \mathbf{A}\mathbf{B}) = \sum_{k=0}^{\infty} \text{Tr} \wedge^k(\mathbf{A}\mathbf{B}) \quad (13)$$

$$= \sum_{k=0}^{\infty} \text{Tr} \wedge^k(\mathbf{A}) \wedge^k(\mathbf{B}) \quad (14)$$

$$= \sum_{k=0}^{\infty} \text{Tr} \wedge^k(\mathbf{B}) \wedge^k(\mathbf{A}) \quad (15)$$

$$= \det(\mathbf{I} + \mathbf{B}\mathbf{A}) \quad (16)$$

where trace class ensures that the sums manipulated are finite. \blacksquare

1.3 Connections to Jackknife

Wu et al. (1986) extend the analogy of the LSE being the mean of the volume sampling distribution

$$\hat{\beta} = \sum_{|S|=k} \det(\mathbf{X}_S^\top \mathbf{X}_S) \hat{\beta}_S = \mathbb{E} \hat{\beta}_S \quad (17)$$

to second-order variance estimators and show (for $\mathbf{y} = \mathbf{X}\beta + \xi$, $\xi \sim \mathcal{N}(0, \sigma^2)$)

$$\mathbb{E}_\xi \sum_{|S|=k} \det(\mathbf{X}_S^\top \mathbf{X}_S) (\hat{\beta}_S - \hat{\beta})(\hat{\beta}_S - \hat{\beta}^\top) = \binom{n-d}{k-d+1} \det(\mathbf{X}^\top \mathbf{X})(\mathbf{X}^\top \mathbf{X})^{-1} \quad (18)$$

In other words, the variance estimator is unbiased for homoscedastic noise.

Bibliography

Dobriban, E. and Y. Sheng

2019. One-shot distributed ridge regression in high dimensions. *arXiv preprint arXiv:1903.09321*.

Wu, C.-F. J. et al.

1986. Jackknife, bootstrap and other resampling methods in regression analysis. *the Annals of Statistics*, 14(4):1261–1295.