

1) **(Chapter 15)** Consider a first-order differentiable convex risk function $P(w)$ with bounded gradients, and introduce the convex optimization problem

$$w^* = \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} P(w) \quad \text{subject to } w \in \mathcal{C}$$

where \mathcal{C} is a convex set. The bounded gradient assumption amounts to

$$\|\nabla_w P(w)\|^2 \leq B, \quad \forall w \in \mathcal{C}$$

for some constant $B \geq 0$. Introduce the projection gradient algorithm

$$\begin{aligned} z_n &= w_{n-1} - \mu \nabla_{w^\top} P(w_{n-1}) \\ w_n &= \mathcal{P}_C[z_n] \end{aligned}$$

where \mathcal{P}_C refers to the projection operator onto \mathcal{C} .

(a) Establish that

$$\frac{1}{N+1} \left(\sum_{m=0}^N (P(w_{m-1}) - P(w^*)) \right) \leq \frac{\|\tilde{w}_{-1}\|^2 + \mu^2 B^2 (N+1)}{2(N+1)\mu}$$

(b) Assume we run the algorithm up to iteration N and return

$$w_N^{\text{best}} = \underset{0 \leq n \leq N}{\operatorname{argmin}} P(w_{n-1})$$

That is, we return the vector with the smallest risk value. Show that if $N = O(1/\epsilon^2)$ then

$$P(w_N^{\text{best}}) \leq P(w^*) + \epsilon$$

(c) Assume we return instead the average weight

$$w_N^{\text{av}} = \frac{1}{N+1} \sum_{n=0}^N w_{n-1}$$

Show again that if $N = O(1/\epsilon^2)$ and $\mu = O(1/\sqrt{N})$ then

$$P(w_N^{\text{av}}) \leq P(w^*) + \epsilon$$

Solution:

(a) Note that

$$\begin{aligned} w_n &= \mathcal{P}_C[w_{n-1} - \mu \nabla_{w^\top} P(w_{n-1})] \\ w^* &= \mathcal{P}_C[w^*] \end{aligned}$$

Using the nonexpansive property (9.70) we get

$$\begin{aligned} \|w^* - w_n\|^2 &\leq \left\| \mathcal{P}_C[w^*] - \mathcal{P}_C[w_{n-1} - \mu \nabla_{w^\top} P(w_{n-1})] \right\|^2 \\ &\leq \|w^* - w_{n-1} + \mu \nabla_{w^\top} P(w_{n-1})\|^2 \end{aligned}$$

so that

$$\begin{aligned}\|\tilde{w}_n\|^2 &\leq \|\tilde{w}_{n-1} + \mu \nabla_{w^\top} P(w_{n-1})\|^2 \\ &\leq \|\tilde{w}_{n-1}\|^2 + \mu^2 \|\nabla_{w^\top} P(w_{n-1})\|^2 + 2\mu \nabla_w P(w_{n-1}) \tilde{w}_{n-1} \\ &\leq \|\tilde{w}_{n-1}\|^2 + \mu^2 B^2 + 2\mu \nabla_w P(w_{n-1}) \tilde{w}_{n-1}\end{aligned}$$

Next, from convexity we have

$$P(w^*) \geq P(w_{n-1}) + \nabla_w P(w_{n-1})(w^* - w_{n-1})$$

That is,

$$\nabla_w P(w_{n-1})(w^* - w_{n-1}) \leq P(w^*) - P(w_{n-1})$$

It follows that

$$\|\tilde{w}_n\|^2 \leq \|\tilde{w}_{n-1}\|^2 + \mu^2 B^2 - 2\mu (P(w_{n-1}) - P(w^*))$$

Iterating gives

$$0 \leq \|\tilde{w}_N\|^2 \leq \|\tilde{w}_{-1}\|^2 + \mu^2 B^2 (N+1) - 2\mu \sum_{m=0}^N (P(w_{m-1}) - P(w^*))$$

which implies that

$$\frac{1}{N+1} \left(\sum_{m=0}^N (P(w_{m-1}) - P(w^*)) \right) \leq \frac{\|\tilde{w}_{-1}\|^2 + \mu^2 B^2 (N+1)}{2(N+1)\mu}$$

(b) For

$$w_N^{\text{best}} = \underset{0 \leq n \leq N}{\operatorname{argmin}} P(w_{n-1})$$

we have

$$P(w_N^{\text{best}}) - P(w^*) \leq \frac{1}{N+1} \left(\sum_{m=0}^N (P(w_{m-1}) - P(w^*)) \right)$$

and, hence,

$$P(w_N^{\text{best}}) - P(w^*) \leq \frac{\|\tilde{w}_{-1}\|^2}{2(N+1)} \frac{1}{\mu} + \frac{\mu B^2}{2}$$

Differentiating the upper bound relative to μ gives $\mu = O(1/\sqrt{N})$ from which

$$P(w_N^{\text{best}}) - P(w^*) \leq O\left(\frac{1}{\sqrt{N}}\right)$$

For the upper bound to be $O(\epsilon)$ we need $N = O(1/\epsilon^2)$.

(c) By convexity we have

$$P(w_N^{\text{av}}) = P\left(\frac{1}{N+1} \sum_{n=0}^N w_{n-1}\right) \leq \frac{1}{N+1} \sum_{n=0}^N P(w_{n-1})$$

so that

$$P(w_N^{\text{av}}) - P(w^*) \leq \frac{\|\tilde{w}_{-1}\|^2 + \mu^2 B^2 (N+1)}{2(N+1)\mu}$$

That is,

$$P(w_N^{\text{av}}) - P(w^*) \leq \frac{\|\tilde{w}_{-1}\|^2}{2(N+1)} \frac{1}{\mu} + \frac{\mu B^2}{2}$$

Differentiating the upper bound relative to μ gives $\mu = O(1/\sqrt{N})$ from which

$$P(w_N^{\text{av}}) - P(w^*) \leq O\left(\frac{1}{\sqrt{N}}\right)$$

For the upper bound to be $O(\epsilon)$ we need $N = O(1/\epsilon^2)$. □

2) **(Chapter 17)** Consider the following variation of Polyak's momentum acceleration method described originally by (17.65), where the second step involves now a combination of \bar{b}_{n-1} and b_n using a smoothing factor β :

$$\begin{aligned} b_n &= \nabla_{w^\top} Q(w_{n-1}; \gamma(n), h_n) \\ \bar{b}_n &= \beta \bar{b}_{n-1} + (1 - \beta) b_n \\ w_n &= w_{n-1} - \mu \bar{b}_n \end{aligned}$$

(a) Assume b_n and \bar{b}_{n-1} are unbiased estimators for the true gradient vectors at w_{n-1} and w_{n-2} , respectively, i.e., $\mathbb{E} b_n = \mathbb{E} \nabla_{w^\top} P(w_{n-1})$ and $\mathbb{E} \bar{b}_{n-1} = \mathbb{E} \nabla_{w^\top} P(w_{n-2})$. Is the updated vector \bar{b}_n an unbiased estimator for $\nabla_{w^\top} P(w_{n-1})$?

(b) Assume we modify the second step as follows:

$$\bar{b}_n = \beta \left[\bar{b}_{n-1} - \mathbb{E} \nabla_{w^\top} P(w_{n-2}) + \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \right] + (1 - \beta) b_n$$

Will \bar{b}_n be now an unbiased estimator for $\nabla_{w^\top} P(w_n)$?

(c) Use a Taylor series expansion to explain that the correction in part (b) can be approximated by the following in terms of the Hessian matrix of the empirical risk evaluated at w_{n-1} :

$$\bar{b}_n = \beta \left[\bar{b}_{n-1} + \mathbb{E} \nabla_w^2 P(w_{n-1})(w_{n-1} - w_{n-2}) \right] + (1 - \beta) b_n$$

(d) What would be a stochastic approximation for the construction in part (c)?

Solution: For more details, the reader may refer to the works by Cutkosky, H. and F. Orabona (2019), “Momentum-based variance reduction in nonconvex SGD,” *Proc. Advances Neural Information Processing Systems*, pp. 1–10, and Tran, H. and A. Cutkosky (2022), “Better SGD using second-order momentum” *Proc. Advances Neural Information Processing Systems*, pp. 3530–3541.

(a) Note that

$$\begin{aligned} \mathbb{E} \bar{b}_n &= \beta \mathbb{E} \bar{b}_{n-1} + (1 - \beta) \mathbb{E} b_n \\ &= \beta \mathbb{E} \nabla_{w^\top} P(w_{n-2}) + (1 - \beta) \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \\ &\neq \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \end{aligned}$$

It follows that the variable \bar{b}_n is not an unbiased estimator for $\nabla_{w^\top} P(w_{n-1})$.

(b) In this case we have

$$\begin{aligned} \mathbb{E} \bar{b}_n &= \beta \left[\mathbb{E} \bar{b}_{n-1} - \mathbb{E} \nabla_{w^\top} P(w_{n-2}) + \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \right] + (1 - \beta) \mathbb{E} b_n \\ &= \beta \left[\mathbb{E} \nabla_{w^\top} P(w_{n-2}) - \mathbb{E} \nabla_{w^\top} P(w_{n-2}) + \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \right] + (1 - \beta) \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \\ &= \mathbb{E} \nabla_{w^\top} P(w_{n-1}) \end{aligned}$$

It follows that \bar{b}_n is now an unbiased estimator for $\nabla_{w^\top} P(w_{n-1})$?

(c) Using the Taylor series expansion we have

$$\nabla_{w^\top} P(w_{n-2}) = \nabla_{w^\top} P(w_{n-1}) + \nabla_w^2 P(w_{n-1})(w_{n-2} - w_{n-1}) + O(\|w_{n-2} - w_{n-1}\|^2)$$

Ignoring the higher order term and substituting into part (b) gives

$$\bar{\mathbf{b}}_n = \beta \left[\bar{\mathbf{b}}_{n-1} + \mathbb{E} \nabla_w^2 P(\mathbf{w}_{n-1})(\mathbf{w}_{n-1} - \mathbf{w}_{n-2}) \right] + (1 - \beta) \mathbf{b}_n$$

(d) A stochastic approximation replaces the Hessian matrix of $P(w)$ by the Hessian matrix of the loss function to get

$$\bar{\mathbf{b}}_n = \beta \left[\bar{\mathbf{b}}_{n-1} + \nabla_w^2 Q(\mathbf{w}_{n-1}; \boldsymbol{\gamma}(n), \mathbf{h}_n)(\mathbf{w}_{n-1} - \mathbf{w}_{n-2}) \right] + (1 - \beta) \mathbf{b}_n$$

□

3) (Chapter 18) Consider a collection $\mathcal{S} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ of N independent random vectors, each with mean $\mathbb{E} \mathbf{x}_n = \bar{\mathbf{x}}_n$ and variance $\sigma_n^2 = \mathbb{E} \|\mathbf{x}_n - \mathbb{E} \mathbf{x}_n\|^2$. Introduce the sample mean

$$\bar{\mathbf{x}} \triangleq \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

We wish to estimate $\bar{\mathbf{x}}$ by considering two constructions based on a mini-batch of size B . The first construction selects B at random *without* replacement from \mathcal{S} , while the second construction selects B samples *with* replacement. We do not assume uniform sampling. Instead, we let p_n denote the normalized *inclusion* probability of \mathbf{x}_n in the mini-batch \mathcal{B} . We denote the samples selected with replacement by $\{\mathbf{x}_b^r\}$ and the samples selected without replacement by $\{\mathbf{x}_b^{wr}\}$. We construct the estimators as follows:

$$\begin{aligned} \hat{\mathbf{x}}^r &\triangleq \frac{1}{B} \sum_{b=1}^B \frac{1}{Np_b} \mathbf{x}_b^r \quad (\text{with replacement}) \\ \hat{\mathbf{x}}^{wr} &\triangleq \frac{1}{B} \sum_{b=1}^B \frac{1}{Np_b} \mathbf{x}_b^{wr} \quad (\text{without replacement}) \end{aligned}$$

Show that both estimators are unbiased with variances

$$\begin{aligned} \mathbb{E} \|\hat{\mathbf{x}}^r - \bar{\mathbf{x}}\|^2 &= \frac{1}{B} \sum_{n=1}^N p_n \left(\frac{1}{N^2 p_n^2} \sigma_n^2 + \left\| \frac{1}{Np_n} \bar{\mathbf{x}}_n - \bar{\mathbf{x}} \right\|^2 \right) \\ \mathbb{E} \|\hat{\mathbf{x}}^{wr} - \bar{\mathbf{x}}\|^2 &\leq \frac{1}{B} \sum_{n=1}^N p_n \left(\frac{1}{N^2 p_n^2} \sigma_n^2 + \left\| \frac{1}{Np_n} \bar{\mathbf{x}}_n - \bar{\mathbf{x}} \right\|^2 \right) \end{aligned}$$

Solution: The argument appears in the work by Rizk, E., S. Vlaski, and A. H. Sayed (2022), “Federated learning under importance sampling,” *IEEE Trans. Signal Process.*, vol. 70, pp. 5381–5396. The derivation extends a result from Brewer, K. R. W. and M. Hanif (1983), *Sampling with Unequal Probabilities*, Springer-Verlag, NY. The presentation is adjusted for the current notation and context.

We begin with the *with-replacement* setting. The randomness of the samples introduces some intri-

cacies that need to be accounted for in the notation. For the mean, we have:

$$\begin{aligned}
\mathbb{E}\widehat{\mathbf{x}}^r &= \frac{1}{B} \sum_{b=1}^B \mathbb{E} \left(\frac{1}{Np_b} \mathbf{x}_b^r \right) \\
&= \frac{1}{B} \sum_{b=1}^B \mathbb{E} \left\{ \mathbb{E} \left\{ \frac{1}{Np_b} \mathbf{x}_b^r \middle| \mathcal{S} \right\} \right\} \\
&= \frac{1}{B} \sum_{b=1}^B \mathbb{E} \left\{ \sum_{n=1}^N p_n \frac{1}{Np_n} \mathbf{x}_n \right\} = \frac{1}{B} \sum_{b=1}^B \bar{x} \\
&= \bar{x}.
\end{aligned}$$

For the variance we find:

$$\begin{aligned}
&\mathbb{E}\|\widehat{\mathbf{x}}^r - \bar{x}\|^2 \\
&= \mathbb{E} \left\| \frac{1}{B} \sum_{b=1}^B \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2 \\
&= \mathbb{E} \left\| \frac{1}{B} \sum_{b=1}^B \left(\frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right) \right\|^2, \\
&= \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\| \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2 \\
&\quad + \frac{1}{B^2} \sum_{b_1 \neq b_2} \mathbb{E} \left\{ \left(\frac{1}{Np_{b_1}} \mathbf{x}_{b_1}^r - \bar{x} \right) \left(\frac{1}{Np_{b_2}} \mathbf{x}_{b_2}^r - \bar{x} \right) \right\} \\
&\stackrel{(a)}{=} \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\| \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2 \\
&\quad + \frac{1}{B^2} \sum_{b_1 \neq b_2} \mathbb{E} \left\{ \frac{1}{Np_{b_1}} \mathbf{x}_{b_1}^r - \bar{x} \right\} \mathbb{E} \left\{ \frac{1}{Np_{b_2}} \mathbf{x}_{b_2}^r - \bar{x} \right\}, \\
&\stackrel{(b)}{=} \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\| \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2
\end{aligned}$$

where (a) is a result of the fact that the elements of \mathcal{S} are independent and \mathbf{x}_b^r is sampled from \mathcal{S} independently, and hence \mathbf{x}_{b_1} and \mathbf{x}_{b_2} are independent. Step (b) follows from:

$$\mathbb{E} \left(\frac{1}{Np_b} \mathbf{x}_b \right) = \mathbb{E} \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) = \bar{x}.$$

Then,

$$\begin{aligned}
& \mathbb{E} \|\hat{\mathbf{x}}^r - \bar{x}\|^2 \\
&= \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\| \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2 \\
&= \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\{ \mathbb{E} \left\| \frac{1}{Np_b} \mathbf{x}_b^r - \bar{x} \right\|^2 \middle| \mathcal{S} \right\} \\
&= \frac{1}{B^2} \sum_{b=1}^B \mathbb{E} \left\{ \sum_{n=1}^N p_n \left\| \frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right\|^2 \right\}, \\
&= \frac{1}{B^2} \sum_{b=1}^B \sum_{n=1}^N p_n \mathbb{E} \left\| \frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right\|^2 \\
&= \frac{1}{B} \sum_{n=1}^N p_n \mathbb{E} \left\| \frac{1}{Np_n} \mathbf{x}_n - \frac{1}{Np_n} \bar{x}_n + \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \\
&= \frac{1}{B} \sum_{n=1}^N p_n \left(\mathbb{E} \left\| \frac{1}{Np_n} \mathbf{x}_n - \frac{1}{Np_n} \bar{x}_n \right\|^2 + \left\| \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \right) \\
&= \frac{1}{B} \sum_{n=1}^N p_n \left(\frac{1}{N^2 p_n^2} \sigma_n^2 + \left\| \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \right).
\end{aligned}$$

We now proceed to study the *without replacement* setting. The fact that the \mathbf{x}_b are sampled from \mathcal{S} without replacement causes pairs $\mathbf{x}_{b_1}, \mathbf{x}_{b_2}$ to no longer be independent. We denote the set of points sampled from \mathcal{S} *without replacement* by \mathcal{B}^{wr} and introduce the activation function by:

$$\mathbb{I}_n \triangleq \begin{cases} 1, & \text{if } \mathbf{x}_n \in \mathcal{B}^{\text{wr}} \\ 0, & \text{if } \mathbf{x}_n \notin \mathcal{B}^{\text{wr}} \end{cases}$$

Then, the estimator $\hat{\mathbf{x}}^{\text{wr}}$ can be written equivalently as:

$$\hat{\mathbf{x}}^{\text{wr}} = \frac{1}{B} \sum_{n=1}^N \mathbb{I}_n \frac{1}{Np_n} \mathbf{x}_n$$

For the mean, we have:

$$\begin{aligned}
\mathbb{E} \hat{\mathbf{x}}^{\text{wr}} &= \frac{1}{B} \sum_{n=1}^N \mathbb{E} \left\{ \mathbb{I}_n \frac{1}{Np_n} \mathbf{x}_n \right\} = \frac{1}{B} \sum_{n=1}^N \mathbb{E} \mathbb{I}_n \times \mathbb{E} \frac{1}{Np_n} \mathbf{x}_n \\
&= \frac{1}{B} \sum_{n=1}^N Bp_n \times \frac{1}{Np_n} \bar{x}_n = \frac{1}{N} \sum_{n=1}^N \bar{x}_n = \bar{x}
\end{aligned}$$

For the variance, we have:

$$\begin{aligned}
\mathbb{E} \|\hat{\mathbf{x}}^{\text{wr}} - \bar{x}\|^2 &= \mathbb{E} \left\| \frac{1}{B} \sum_{n=1}^N \mathbb{I}_n \left(\frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right) \right\|^2 \\
&= \frac{1}{B^2} \sum_{n=1}^N \mathbb{E} \left\| \mathbb{I}_n \left(\frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right) \right\|^2 \\
&\quad + \frac{1}{B^2} \sum_{n_1 \neq n_2} \mathbb{E} \left\{ \mathbb{I}_{n_1} \left(\frac{1}{Np_{n_1}} \mathbf{x}_{n_1} - \bar{x} \right) \mathbb{I}_{n_2} \right. \\
&\quad \left. \times \left(\frac{1}{Np_{n_2}} \mathbf{x}_{n_2} - \bar{x} \right) \right\}
\end{aligned}$$

We begin with:

$$\begin{aligned}
&\mathbb{E} \left\| \mathbb{I}_n \left(\frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right) \right\|^2 \\
&= \mathbb{E} \left\{ \left\| \mathbb{I}_n \left(\frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right) \right\|^2 \middle| \mathbb{I}_n = 1 \right\} \times \mathbb{P}(\mathbb{I}_n = 1) \\
&\quad + \mathbb{E} \left\{ \left\| \mathbb{I}_n \left(\frac{1}{Np_n} \mathbf{x}_n - \bar{x} \right) \right\|^2 \middle| \mathbb{I}_n = 0 \right\} \times \mathbb{P}(\mathbb{I}_n = 0) \\
&= Bp_n \left(\mathbb{E} \left\| \frac{1}{Np_n} \mathbf{x}_n - \frac{1}{Np_n} \bar{x}_n + \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \right) \\
&= Bp_n \left(\frac{1}{N^2 p_n^2} \mathbb{E} \|\mathbf{x}_n - \bar{x}_n\|^2 + \left\| \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \right) \\
&= Bp_n \left(\frac{1}{N^2 p_n^2} \sigma_n^2 + \left\| \frac{1}{Np_n} \bar{x}_n - \bar{x} \right\|^2 \right)
\end{aligned}$$

For the cross-term we have:

$$\begin{aligned}
&\mathbb{E} \left\{ \mathbb{I}_{n_1} \left(\frac{1}{Np_{n_1}} \mathbf{x}_{n_1} - \bar{x} \right) \mathbb{I}_{n_2} \left(\frac{1}{Np_{n_2}} \mathbf{x}_{n_2} - \bar{x} \right) \right\} \\
&= \mathbb{E} \left\{ \left(\frac{1}{Np_{n_1}} \mathbf{x}_{n_1} - \bar{x} \right) \left(\frac{1}{Np_{n_2}} \mathbf{x}_{n_2} - \bar{x} \right) \middle| \mathbb{I}_{n_1} = 1, \mathbb{I}_{n_2} = 1 \right\} \\
&\quad \times \mathbb{P}(\mathbb{I}_{n_1} = 1, \mathbb{I}_{n_2} = 1) \\
&= \mathbb{P}(\mathbb{I}_{n_2} = 1, \mathbb{I}_{n_1} = 1) \left(\frac{1}{Np_{n_1}} \mathbb{E} \mathbf{x}_{n_1} - \bar{x} \right) \left(\frac{1}{Np_{n_2}} \mathbb{E} \mathbf{x}_{n_2} - \bar{x} \right) \\
&= \mathbb{P}(\mathbb{I}_{n_2} = 1, \mathbb{I}_{n_1} = 1) \left(\frac{1}{Np_{n_1}} \bar{x}_{n_1} - \bar{x} \right) \left(\frac{1}{Np_{n_2}} \bar{x}_{n_2} - \bar{x} \right)
\end{aligned}$$

We then get the desired result:

$$\begin{aligned}
& \mathbb{E} \|\hat{\mathbf{x}}^{\text{wr}} - \bar{\mathbf{x}}\|^2 \\
&= \frac{1}{B} \sum_{n=1}^N p_n \left(\frac{1}{N^2 p_n^2} \sigma_n^2 + \left\| \frac{1}{N p_n} \bar{\mathbf{x}}_n - \bar{\mathbf{x}} \right\|^2 \right) \\
&\quad + \frac{1}{B^2} \sum_{n_1 \neq n_2} \mathbb{P} \{ \mathbb{I}_{n_2} = 1, \mathbb{I}_{n_1} = 1 \} \left(\frac{1}{N p_{n_1}} \bar{\mathbf{x}}_{n_1} - \bar{\mathbf{x}} \right) \\
&\quad \times \left(\frac{1}{N p_{n_2}} \bar{\mathbf{x}}_{n_2} - \bar{\mathbf{x}} \right).
\end{aligned}$$

□

4) **(Chapter 19)** Consider an empirical risk minimization problem of the following form where $P(w)$ is ν -strongly convex and has δ -Lipschitz gradients:

$$w^* \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ P(w) \triangleq \frac{1}{N} \sum_{m=0}^{N-1} Q(w; \gamma(m), h_m) \right\}$$

Consider the average regret from the body of the chapter defined as

$$\mathbb{E} \mathcal{R}(N) \triangleq \frac{1}{N} \sum_{n=0}^{N-1} \left(\mathbb{E} P(\mathbf{w}_{n-1}) - P(w^*) \right)$$

The empirical risk is minimized by using the classical stochastic gradient algorithm with step-size μ and where the gradient noise satisfies conditions (19.13a)–(19.13b).

(a) Show that an approximate upper bound for $\mathbb{E} \mathcal{R}(N)$ is given by

$$\mathbb{E} \mathcal{R}(N) \leq \frac{\delta}{4N\mu\nu} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu\sigma_g^2\delta}{4\nu}$$

(b) Minimize the upper bound with respect to μ and conclude that $\mathbb{E} \mathcal{R}(N) \leq O(1/\sqrt{N})$.

Solution:

(a) We know that

$$\begin{aligned}
0 &\leq \mathbb{E} P(\mathbf{w}_{n-1}) - P(w^*) \leq \frac{\delta}{2} \|\tilde{\mathbf{w}}_{n-1}\|^2 \\
\|\mathbb{E} \tilde{\mathbf{w}}_{n-1}\|^2 &\leq \lambda^n \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu^2 \sigma_g^2}{1-\lambda} \approx \lambda^n \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu\sigma_g^2}{2\nu}
\end{aligned}$$

where we approximated $1 - \lambda \approx 2\mu\nu$. It follows that

$$\mathbb{E} P(\mathbf{w}_{n-1}) - P(w^*) \leq \frac{\delta}{2} \lambda^n \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu\sigma_g^2\delta}{4\nu}$$

and, hence,

$$\begin{aligned}
\mathbb{E} \mathcal{R}(N) &\leq \frac{1}{N} \sum_{n=0}^{N-1} \left(\frac{\delta}{2} \lambda^n \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu \sigma_g^2 \delta}{4\nu} \right) \\
&\leq \frac{\delta}{2N} \frac{1}{1-\lambda} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu \sigma_g^2 \delta}{4\nu} \\
&\approx \frac{\delta}{2N} \frac{1}{2\mu\nu} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu \sigma_g^2 \delta}{4\nu} \\
&= \frac{\delta}{4N\mu\nu} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\mu \sigma_g^2 \delta}{4\nu}
\end{aligned}$$

(b) Differentiating the upper bound relative to μ and setting the derivative to 0 gives:

$$-\frac{\delta}{4N\nu} \frac{1}{\mu^2} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{\sigma_g^2 \delta}{4\nu} = 0$$

for which

$$\mu^o = \frac{\left(\mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 \right)^{1/2}}{\sigma_g \sqrt{N}}$$

Substituting into the upper bound for $\mathbb{E} \mathcal{R}(N)$ gives

$$\mathbb{E} \mathcal{R}(N) \leq \frac{\delta \sigma_g \left(\mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 \right)^{1/2}}{2\nu} \times \frac{1}{\sqrt{N}} = O(1/\sqrt{N})$$

□

5) (Chapter 27) Consider the noisy observation $\mathbf{y} = \mathbf{x} + \mathbf{v}$ where all variables are M -dimensional and \mathbf{v} is Gaussian noise with zero mean and covariance matrix R_v , i.e., $\mathbf{v} \sim \mathcal{N}(0, R_v)$. The noise and \mathbf{x} are independent of each other and R_v is positive definite.

(a) Let $f_{\mathbf{x}}(x)$ denote the pdf of \mathbf{x} . Argue that the pdf of \mathbf{y} is given by

$$f_{\mathbf{y}}(y) \triangleq \int_{x \in \mathcal{X}} \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp \left\{ -\frac{1}{2} (y - x)^\top R_v^{-1} (y - x) \right\} f_{\mathbf{x}}(x) dx$$

Argue further that the conditional pdf of \mathbf{x} given \mathbf{y} is

$$f_{\mathbf{x}|\mathbf{y}}(x|y) = \frac{f_{\mathbf{x}}(x)}{f_{\mathbf{y}}(y)} \times \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp \left\{ -\frac{1}{2} (y - x)^\top R_v^{-1} (y - x) \right\}$$

(b) Verify that the gradient of $f_{\mathbf{y}}(y)$ relative to y satisfies

$$\frac{1}{f_{\mathbf{y}}(y)} \nabla_{y^\top} f_{\mathbf{y}}(y) = -R_v^{-1} (y - \hat{x})$$

where $\hat{x} = \mathbb{E}(\mathbf{x}|\mathbf{y} = y)$.

(c) Conclude that the least mean-squares estimator satisfies the following relation in terms of the pdf of the observation:

$$\mathbb{E}(\mathbf{x}|\mathbf{y}) = \mathbf{y} + R_v \nabla_{y^\top} \ln f_{\mathbf{y}}(y)$$

Solution:

(a) From Bayes rule

$$f_{\mathbf{x},\mathbf{y}}(x, y) = f_{\mathbf{x}}(x) f_{\mathbf{y}|\mathbf{x}}(y|x) = f_{\mathbf{y}}(y) f_{\mathbf{x}|\mathbf{y}}(x|y)$$

Moreover, we have that

$$f_{\mathbf{y}|\mathbf{x}}(y|x) = f_{\mathbf{v}}(y-x) = \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\}$$

Therefore,

$$\begin{aligned} f_{\mathbf{x},\mathbf{y}}(x,y) &= f_{\mathbf{x}}(x) f_{\mathbf{y}|\mathbf{x}}(y|x) \\ &= f_{\mathbf{x}}(x) \times \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} \end{aligned}$$

from which we conclude by marginalizing over x that

$$f_{\mathbf{y}}(y) = \int_{x \in \mathcal{X}} \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} f_{\mathbf{x}}(x) dx$$

Likewise, from Bayes rule again

$$\begin{aligned} f_{\mathbf{x}|\mathbf{y}}(x|y) &= \frac{f_{\mathbf{x}}(x) f_{\mathbf{y}|\mathbf{x}}(y|x)}{f_{\mathbf{y}}(y)} \\ &= \frac{f_{\mathbf{x}}(x)}{f_{\mathbf{y}}(y)} \times \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} \end{aligned}$$

(b) Differentiating $f_{\mathbf{y}}(y)$ relative to y and switching differentiation with integration gives

$$\begin{aligned} &\nabla_{y^\top} f_{\mathbf{y}}(y) \\ &= \nabla_{y^\top} \left\{ \int_{x \in \mathcal{X}} \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} f_{\mathbf{x}}(x) dx \right\} \\ &= \int_{x \in \mathcal{X}} \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \nabla_{y^\top} \left\{ \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} \right\} f_{\mathbf{x}}(x) dx \\ &= -\frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \int_{x \in \mathcal{X}} R_v^{-1}(y-x) \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} f_{\mathbf{x}}(x) dx \\ &= -R_v^{-1} \times \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \int_{x \in \mathcal{X}} y \times \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} f_{\mathbf{x}}(x) dx + \\ &\quad R_v^{-1} \times \frac{1}{\sqrt{(2\pi)^M}} \frac{1}{\sqrt{\det R_v}} \int_{x \in \mathcal{X}} x \times \exp\left\{-\frac{1}{2}(y-x)^\top R_v^{-1}(y-x)\right\} f_{\mathbf{x}}(x) dx \\ &= -R_v^{-1} \times \int_{x \in \mathcal{X}} y \times f_{\mathbf{y}|\mathbf{x}}(y|x) f_{\mathbf{x}}(x) dx + R_v^{-1} \times f_{\mathbf{y}}(y) \int_{x \in \mathcal{X}} x \times f_{\mathbf{x}|\mathbf{y}}(x|y) dx \\ &= -R_v^{-1} \times \int_{x \in \mathcal{X}} y \times f_{\mathbf{x},\mathbf{y}}(x,y) dx + R_v^{-1} \times f_{\mathbf{y}}(y) \times \mathbb{E}(\mathbf{x}|\mathbf{y} = y) \\ &= -R_v^{-1} \times f_{\mathbf{y}}(y) \times \int_{x \in \mathcal{X}} y \times f_{\mathbf{x}|\mathbf{y}}(x|y) dx + R_v^{-1} \times f_{\mathbf{y}}(y) \mathbb{E}(\mathbf{x}|\mathbf{y} = y) \\ &= -R_v^{-1} \times f_{\mathbf{y}}(y) \times y \times \underbrace{\int_{x \in \mathcal{X}} f_{\mathbf{x}|\mathbf{y}}(x|y) dx}_{=1} + R_v^{-1} \times f_{\mathbf{y}}(y) \mathbb{E}(\mathbf{x}|\mathbf{y} = y) \\ &= R_v^{-1} \times f_{\mathbf{y}}(y) \times (\mathbb{E}(\mathbf{x}|\mathbf{y} = y) - y) \end{aligned}$$

Therefore,

$$\frac{1}{f_{\mathbf{y}}(y)} \nabla_{y^\top} f_{\mathbf{y}}(y) = R_v^{-1} \times f_{\mathbf{y}}(y) \times (\mathbb{E}(\mathbf{x}|\mathbf{y} = y) - y)$$

where $\hat{x} = \mathbb{E}(\mathbf{x}|\mathbf{y} = y)$.

(c) Rearranging terms and noting that

$$\nabla_{y^\top} \ln f_{\mathbf{y}}(y) = \frac{1}{f_{\mathbf{y}}(y)} \nabla_{y^\top} f_{\mathbf{y}}(y)$$

we conclude that

$$\mathbb{E}(\mathbf{x}|\mathbf{y}) = \mathbf{y} + R_v \nabla_{y^\top} \ln f_{\mathbf{y}}(y)$$

□