# Supplementary material: Following the Perturbed Leader for Online Structured Learning

## 7. Appendix

### 7.1. Additional proofs

#### 7.1.1. MAXIMA OF NORMAL RANDOM VARIABLES

**Lemma 9.** *Suppose the conditions of Theorem 1 hold, then*

$$\mathbb{E}\left[\max_{x\in\mathcal{X}}\langle x,\gamma\rangle\right] \le \sqrt{2k\log|\mathcal{X}|}$$

*Proof.* First, we upper bound the expectation by

$$\mathbb{E}_{\gamma}\left[\max_{x\in\mathcal{X}}\langle x,\gamma\rangle\right] \le \inf_{s>0}\frac{1}{s}\log\left(\sum_{x\in\mathcal{X}}\mathbb{E}[\exp(s\langle x,\gamma\rangle)]\right)$$

Notice that $\langle x,\gamma\rangle$ is a normal random variable with mean $0$ and variance $\|x\|^2 \le k$. As such,

$$\mathbb{E}[\exp(s\langle x,\gamma\rangle)] = \exp\left(\frac{s^2\|x\|^2}{2}\right) \le \exp\left(\frac{ks^2}{2}\right)$$

Then,

$$\mathbb{E}_{\gamma}\left[\max_{x\in\mathcal{X}}\langle x,\gamma\rangle\right] \le \inf_{s>0}\frac{1}{s}\log\left(|\mathcal{X}|\exp\left(\frac{ks^2}{2}\right)\right)$$
$$= \inf_{s>0}\left\{\frac{\log|\mathcal{X}|}{s} + \frac{ks}{2}\right\}$$
$$= \sqrt{2k\log|\mathcal{X}|}$$

$\square$

#### 7.1.2. BOUNDING THE HESSIAN

**Lemma 10.** *Suppose that the conditions of Theorem 1 hold. Let $H$ denote the Hessian of $\Phi_\eta$ at an arbitrary $\theta$. Fix some $j\in[d]$. Then,*

$$\sum_{i=1}^{d}|H_{i,j}| \le \frac{k}{\eta}\sum_{x\in\mathcal{X}}|\mathbb{E}\left[\gamma_j\mathbb{1}[\hat{x}=x]\right]|$$

*Proof.* Recall the definition of the Hessian:

$$H_{i,j} = \frac{1}{\eta}\mathbb{E}\left[\hat{x}(\tilde{\theta}_t+\eta\gamma)_i\gamma_j\right]$$

Let us abbreviate $\hat{x}(\theta+\eta\gamma)$ as $\hat{x}$. Then,

$$\eta\sum_{i=1}^{d}|H_{i,j}| = \eta\sum_{i:H_{i,j}>0}H_{i,j} - \eta\sum_{i:H_{i,j}\le0}H_{i,j}$$

$$= \mathbb{E}\left[\left(\sum_{i:H_{i,j}>0}\hat{x}_i - \sum_{i:H_{i,j}\le0}\hat{x}_i\right)\gamma_j\right]$$

$$= \sum_{x\in\mathcal{X}}\mathbb{E}\left[\left(\sum_{i:H_{i,j}>0}\hat{x}_i - \sum_{i:H_{i,j}\le0}\hat{x}_i\right)\gamma_j\mathbb{1}_{[\hat{x}=x]}\right]$$

$$= \sum_{x\in\mathcal{X}}\left(\sum_{i:H_{i,j}>0}x_i - \sum_{i:H_{i,j}\le0}x_i\right)\mathbb{E}\left[\gamma_j\mathbb{1}_{[\hat{x}=x]}\right]$$

$$\le k\sum_{x\in\mathcal{X}}|\mathbb{E}\left[\gamma_j\mathbb{1}[\hat{x}=x]\right]|$$

as $\left|\sum_{i:H_{i,j}>0}x_i - \sum_{i:H_{i,j}\le0}x_i\right| \le k$ by assumption.

$\square$

#### 7.1.3. BOUNDING THE HESSIAN FOR THE $k$-SETS PROBLEM

*Proof of lemma 3.* Let $H = \nabla^2\Phi_\eta(\tilde{\theta})$. We have that,

$$H_{i,j} = \frac{1}{\eta}\mathbb{E}\left[\hat{x}(\tilde{\theta}+\eta\gamma)_i\gamma_j\right]$$

with $\hat{x}(z) \in \arg\min_{x\in\mathcal{X}}\langle x,z\rangle$ (Abernethy et al., 2014, Lemma 7). We shall abbreviate $\hat{x}$ for $\hat{x}(\tilde{\theta}+\eta\gamma)$ in the remainder of the proof.

First, notice that

$$\sum_{i,j}H_{i,j} = \frac{1}{\eta}\sum_{i,j}\mathbb{E}[\gamma_j\hat{x}_i] = \frac{k}{\eta}\sum_{j=1}^{d}\mathbb{E}[\gamma_j] = 0$$

Secondly, we argue about the sign of $\mathbb{E}[\gamma_j\hat{x}_i]$. We claim that it is negative if $i = j$ and positive otherwise. To see that, notice that $\gamma_j$ is a symmetric random variable, so that for each $\alpha > 0$ the density of $\gamma_j$ at $\alpha$ and at $-\alpha$ is the same. If $i \ne j$, the event $\hat{x}_i = 1$ is more probable if $\gamma_j = \alpha$ than when $\gamma_j = -\alpha$. If $i = j$ then the opposite is true.

We have,

$$\sum_{i,j} H_{i,j} = \sum_{i,j:H_{i,j}\geq 0} H_{i,j} - \sum_{i,j:H_{i,j}<0} H_{i,j}$$

$$= -2 \sum_{i,j:H_{i,j}<0} H_{i,j}$$

$$= -2\operatorname{Tr}(H)$$

The rest of the proof follows that of lemma 2.

$\square$

### 7.1.4. TECHNICAL LEMMA

**Lemma 11.** *We have,*

$$\max\left\{\min\left\{\frac{Td}{16}, \frac{d\eta\sqrt{2}}{32}\right\}, \frac{Td}{16}\operatorname{erf}\left(-\frac{\sqrt{d}}{4\eta}\right)\right\}$$

$$\geq \min\left\{0.02Td, 0.05d^{5/4}\sqrt{T}\right\}$$

*Proof.* We get,

$$\max\left\{\min\left\{\frac{Td}{16}, \frac{d\eta\sqrt{2}}{32}\right\}, \frac{Td}{16}\operatorname{erf}\left(\frac{\sqrt{d}}{4\eta}\right)\right\}$$

$$\geq \min\left\{\frac{Td}{16},\right.$$

$$\left.\max\left\{\frac{d\eta\sqrt{2}}{32}, \frac{Td}{16}\operatorname{erf}\left(\frac{\sqrt{d}}{4\eta}\right)\right\}\right\}$$

Notice that erf is nondecreasing and concave on $\mathbb{R}_+$. Then,

$$\inf_{\eta>0} \max\left\{\frac{d\eta\sqrt{2}}{32}, \frac{Td}{16}\operatorname{erf}\left(\frac{\sqrt{d}}{4\eta}\right)\right\}$$

$$\geq \min\left\{\inf_{\eta<\sqrt{d}/4} \frac{Td}{16}\operatorname{erf}\left(\frac{\sqrt{d}}{4\eta}\right),\right.$$

$$\left.\inf_{\eta\geq\sqrt{d}/4} \max\left\{\frac{d\eta\sqrt{2}}{32}, \frac{Td}{16}\operatorname{erf}\left(\frac{\sqrt{d}}{4\eta}\right)\right\}\right\}$$

$$\geq \min\left\{\frac{Td}{16}\operatorname{erf}(1),\right.$$

$$\left.\inf_{\eta\geq\sqrt{d}/4} \max\left\{\frac{d\eta\sqrt{2}}{32}, \frac{Td}{16}\frac{\sqrt{d}}{4\eta}\operatorname{erf}(1)\right\}\right\}$$

$$\geq \min\left\{\frac{Td}{16}\operatorname{erf}(1), \sqrt{\frac{d\sqrt{2}}{32}\frac{Td}{16}\frac{\sqrt{d}}{4}\operatorname{erf}(1)}\right\}$$

$$\geq \min\left\{0.05Td, 0.02d^{5/4}\sqrt{T}\right\}$$

as required.

$\square$

## 7.2. Lipschitz property of certain distributions

### 7.2.1. UNIFORM OVER THE CUBE

Remember that we had required the marginals to have a variance of 1. Therefore WLOG we will take the cube to be $C = [0, 1/\sqrt{3}]^d$. Then,

$$\operatorname{TV}(P,Q) = \sup_A \left|\Pr_P[A] - \Pr_Q[A]\right|$$

$$= \sup_A \left|\frac{1}{\operatorname{Vol}(C)}\int_{x\in A} \mathbb{1}_{[x\in C+\{\mu_P\}]} - \mathbb{1}_{[x\in C+\{\mu_Q\}]}\right|$$

$$\leq \sup_A \frac{1}{\operatorname{Vol}(C)}\int_{x\in A} \left|\mathbb{1}_{[x\in C+\{\mu_P\}]} - \mathbb{1}_{[x\in C+\{\mu_Q\}]}\right|$$

$$\leq \frac{1}{\operatorname{Vol}(C)}\int_{x\in\mathbb{R}^d} \left|\mathbb{1}_{[x\in C+\{\mu_P\}]} - \mathbb{1}_{[x\in C+\{\mu_Q\}]}\right|$$

$$= \frac{\operatorname{Vol}((C+\{\mu_P\}) \triangle (C+\{\mu_Q\}))}{\operatorname{Vol}(C)}$$

$$\leq \frac{2(1/\sqrt{3})^{d-1}\|\mu_P - \mu_Q\|_1}{(1/\sqrt{3})^d} = 2\sqrt{3}\|\mu_P - \mu_Q\|_1$$

so that $L = 2\sqrt{3}$.

We now explain the above bound. Suppose that $C + \{\mu_P\}$ and $C + \{\mu_Q\}$ do not intersect. Then we must have $\|\mu_P - \mu_Q\|_\infty > 1/\sqrt{3}$.

$$\operatorname{Vol}((C+\{\mu_P\}) \triangle (C+\{\mu_Q\})) = 2\left(\frac{1}{\sqrt{3}}\right)^d$$

$$< 2\left(\frac{1}{\sqrt{3}}\right)^{d-1}\|\mu_P - \mu_Q\|_\infty$$

$$\leq 2\left(\frac{1}{\sqrt{3}}\right)^{d-1}\|\mu_P - \mu_Q\|_1$$

If $C + \{\mu_P\}$ and $C + \{\mu_Q\}$ do intersect, then $\|\mu_P - \mu_Q\|_\infty \leq 1/\sqrt{3}$ and we have

$$\operatorname{Vol}((C+\{\mu_P\}) \cap (C+\{\mu_Q\}))$$

$$= \prod_{i=1}^d \left(\frac{1}{\sqrt{3}} - |\mu_{P,i} - \mu_{Q,i}|\right)$$

so that

$$\text{Vol}((C + \{\mu_P\}) \triangle (C + \{\mu_Q\}))$$
$$= \text{Vol}(C + \{\mu_P\}) + \text{Vol}(C + \{\mu_Q\})$$
$$\quad - 2\text{Vol}((C + \{\mu_P\}) \cap (C + \{\mu_Q\}))$$
$$= 2\left(\left(\frac{1}{\sqrt{3}}\right)^d - \prod_{i=1}^{d}\left(\frac{1}{\sqrt{3}} - |\mu_{P,i} - \mu_{Q,i}|\right)\right)$$
$$= 2\left(\frac{1}{\sqrt{3}}\right)^d \left(1 - \prod_{i=1}^{d}\left(1 - \sqrt{3}|\mu_{P,i} - \mu_{Q,i}|\right)\right)$$
$$\leq 2\left(\frac{1}{\sqrt{3}}\right)^d \sqrt{3}\|\mu_P - \mu_Q\|_1$$
$$= 2\left(\frac{1}{\sqrt{3}}\right)^{d-1} \|\mu_P - \mu_Q\|_1$$

### 7.2.2. LAPLACE AND NEGATIVE EXPONENTIAL

We will show that for the Laplace distribution we have $L = \sqrt{2}$. For the exponential distribution the proof is similar except with $L = 1$. Once again, recall that we had required the marginals to have a variance of 1, and therefore the PDF of the Laplace distribution is $\exp(-\sqrt{2}|x - \mu|)/\sqrt{2}$. In this case,

We want to bound

$$\text{TV}(P, Q) = \sup_A \left|\Pr_P[A] - \Pr_Q[A]\right|$$
$$= \sup_A \left|\int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1)\right.$$
$$\left. - \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_Q\|_1)\right|$$

We have,

$$\int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1)$$
$$\quad - \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_Q\|_1)$$
$$= \int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1)$$
$$\quad \cdot \left(1 - \exp(\sqrt{2}\|x - \mu_P\|_1 - \sqrt{2}\|x - \mu_Q\|_1)\right)$$
$$\leq \int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1)$$
$$\quad \cdot \left(\sqrt{2}\|x - \mu_Q\|_1 - \sqrt{2}\|x - \mu_P\|_1\right)$$
$$\leq \int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1) \left(\sqrt{2}\|\mu_P - \mu_Q\|_1\right)$$
$$\leq \sqrt{2}\|\mu_P - \mu_Q\|_1$$

Similarly, one can bound

$$\int_A \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_Q\|_1)$$
$$\quad - \frac{1}{\sqrt{2}} \exp(-\sqrt{2}\|x - \mu_P\|_1) \leq \sqrt{2}\|\mu_P - \mu_Q\|_1$$

and thus

$$\text{TV}(P, Q) \leq \sqrt{2}\|\mu_P - \mu_Q\|_1$$

as claimed.

## References

Abernethy, Jacob, Lee, Chansoo, Sinha, Abhinav, and Tewari, Ambuj. Online linear optimization via smoothing. *The Journal of Machine Learning Research*, 35: 807–823, 2014.