

## Interaction protocol

For  $t = 1, \dots, T$ 

- learner picks  $w_t \in \mathbb{R}^d$ ,  $\|w_t\|_2 = 1$ , in randomized fashion
- environment picks  $L_t \in \mathbb{R}^{d \times d}$ ,  $\|L_t\|_2 \leq 1$
- learner observes loss  $\langle w_t w_t^T, L_t \rangle = \text{tr}(w_t w_t^T L_t)$

Note: -  $L_t$  is loss matrix  
 - allowed to depend on entire history (except  $w_t$ ).  
 - thus, adaptive adversaries allowed.

Goal: minimize regret,

$$R_T = \max_{u: \|u\|_2=1} \sum_{t=1}^T \mathbb{E} [\langle w_t w_t^T - u u^T, L_t \rangle]$$

↳ over randomization of learner

Algorithm Online mirror descentParam:  $\eta > 0$  (learning rate),  $\gamma \in [0, 1]$  (exploration rate)Init:  $W_1 = \frac{1}{d} I_d$ for  $t = 1, \dots, T$ 

$$W_t \stackrel{\text{EVD}}{=} \sum_{i=1}^d \mu_i u_i u_i^T$$

$$\underline{\lambda} = (1-\gamma) \underline{\mu} + \gamma (1/d, \dots, 1/d) \rightarrow \text{controls exploration, exploitation trade-off}$$

$$\tilde{L}_t = \text{sample}(\underline{\lambda}, \{u_i\}_{i=1}^d)$$

↳ 2 possible sub-routines

$$W_{t+1} = (W_t^{-1} + \eta \tilde{L}_t + \beta I)^{-1} \text{ with } \beta \text{ s.t. } \text{tr}(W_{t+1}) = 1$$

$$W_{t+1} = (W_t^{-1} + \eta \tilde{L}_t + \beta I)^{-1} \text{ with } \beta \text{ s.t. } \text{tr}(W_{t+1}) = 1$$

Note:  $W \triangleq (1-\gamma)W_t + \frac{\gamma}{d}I_d$  here

### A. Dense sampling

Sample  $(\lambda, \{u_i\}_{i=1}^d)$

$$B \sim \text{bern}(1/2)$$

to est. diag. elements

if  $B=1$   
 $I \sim \lambda, w_t \leftarrow u_I$

- sample one eig vect s.t.  $P(I=i) = \lambda_i$   
 $E[w_t w_t^T] = E\left[\sum_{i=1}^d \mathbb{1}_{\{I=i\}} u_i u_i^T\right] = \sum \lambda_i u_i u_i^T = W$

to estimate off-diag elements

else  
 $S \in \{-1, 1\}^d$  iid uniformly  
 $w_t \leftarrow \sum_i S_i \sqrt{\lambda_i} u_i$

$$E[w_t w_t^T] = E_S\left[\sum_{i,j} S_i S_j \sqrt{\lambda_i \lambda_j} u_i u_j^T\right] = \sum_{i,j} \delta_{ij} \sqrt{\lambda_i \lambda_j} u_i u_j^T = W$$

play  $w_t$ , observe  $\langle w_t w_t^T, L_t \rangle = d_t$

if  $B=1$

$$\tilde{L}_t \leftarrow 2 d_t W_t^{-1/2} w_t w_t^T W_t^{-1/2}$$

else

$$\tilde{L}_t \leftarrow d_t (W_t^{-1} w_t w_t^T W_t^{-1} - W_t^{-1})$$

return  $\tilde{L}_t$

$$\langle u_i u_i^T, L_t \rangle = \text{tr}(u_i u_i^T L_t) = u_i^T L_t u_i$$

- intuitively, observe we can estimate  $\approx u_i^T L_t u_i$

lemmal  $\rightarrow E[\tilde{L}_t] = L_t$

### Bounding regret

$$R_T = \max_{u: \|u\| \leq 1} \sum_{t=1}^T E[\langle w_t w_t^T - u u^T, L_t \rangle]$$

Theorem 3: let  $\eta \leq \frac{1}{2d}$ ,  $\gamma=0$ . Then,

$$R_T \leq d \frac{\log T}{\eta} + \eta (d^2 + 1) \sum_{t=1}^T E[d_t^2] + 2$$

not so necessary for first read I think

Corollary: let  $\eta = \min\left\{\sqrt{\frac{\log T}{dT}}, \frac{1}{2d}\right\}$ ,  $\gamma=0$

$$R_T \leq O\left(d^{3/2} \sqrt{T \log T}\right)$$

$$R_T \leq O(d^{3/2} \sqrt{T \log T})$$

- if assume that  $L_t \succeq 0$ ,  $\min_{\|u\|_2=1} \sum_t \text{tr}(u u^T L_t) \leq \bar{L}_T^*$   $\hookrightarrow$  approx. says if

then with  $\eta = \min \left\{ \sqrt{\frac{\log T}{d \bar{L}_T^*}}, \frac{1}{4d^2} \right\}$

$\forall t, \|L_t\| = O_T(1)$   
then it gives better bound

$$R_T = O\left(d^{3/2} \sqrt{\bar{L}_T^* \log T} + d^3 \log T\right)$$

### B. Sparse sampling

sample  $(\lambda, \{u_i\}_{i=1}^d)$

draw  $I, J \sim \lambda$

if  $I=J$

$$w_t \leftarrow u_I$$

else

$S \in \{-1, 1\}$  uniformly

$$w_t \leftarrow \frac{1}{\sqrt{2}} (u_I + S u_J)$$

play  $w_t$ , observe  $d_t = \langle w_t w_t^T, L_t \rangle$

if  $I=J$

$$\tilde{L}_t = \frac{d_t}{\lambda_I} u_I u_I^T$$

else

$$\tilde{L}_t = \frac{S d_t}{2\lambda_I \lambda_J} (u_I u_J^T + u_J u_I^T)$$

Same as in dense case,  $E[w_t w_t^T] = W$

little more work, but can show

$$E[w_t w_t^T] = W$$

lemma 2 shows

$$\text{that } E[\tilde{L}_t] = L_t$$

### Bounding regret

Thm 6:  $\eta \leq 1/2d$ ,  $\gamma = \eta d$ .

$$R_T \leq \frac{d \log T}{\eta} + 2hd + 2 + 8hd \sum E[\|L_t\|_F^2]$$

... if  $\eta = \min \left\{ \sqrt{\frac{\log T}{dT}}, \frac{1}{2d} \right\}$

Cor: let  $\frac{1}{T} \sum E[\|L_t\|_F^2] \leq r$ . if  $\eta = \min\left\{\sqrt{\frac{\log T}{rT}}, \frac{1}{2d}\right\}$

$\gamma = nd$ . then

$$R_T = O(d\sqrt{rT \log T})$$

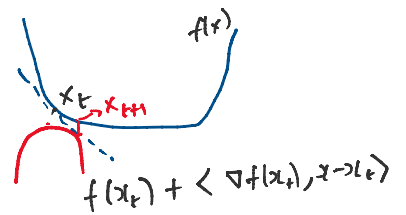
note: since  $\|L_t\|_2 \leq 1$ ,  $\|L_t\|_F \leq \text{rank}(L_t)$

$\Rightarrow$  for  $L_t = x_t x_t^T$  (online PCA),  $R_T \leq O(d\sqrt{rT \log T})$

Understanding the algorithm (Online mirror descent)

(aside) mirror descent! (from Yu's slides)

proximal viewpoint of projected G.D



$$x_{t+1} = \arg \min_{x \in \mathcal{C}} \left\{ \underbrace{f(x_t) + \langle \nabla f(x_t), x - x_t \rangle}_{\text{linear approx}} + \underbrace{\frac{1}{2\eta_t} \|x - x_t\|_2^2}_{\text{prox}} \right\}$$

- the proximal term is based on geometry of problem

- eg. if  $\mathcal{C}$  is probability simplex, better idea to use

- KL divergence, T.V. distance, ...

mirror descent

$$x_{t+1} = \arg \min_{x \in \mathcal{C}} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \underbrace{\frac{1}{\eta_t} D_\Psi(x, x_t)}_{\text{Bregman divergence}} \right\}$$

$$D_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle$$

essentially the difference in first order approximation!

- key point: choice of regularization and associated Bregman is crucial

Back to Bandit PCA:

- Consider differentiable, convex  $R: S \rightarrow \mathbb{R}$ ,  $S \rightarrow d \times d$  symmetric, PSD

- associated  $D_R: S \times S \rightarrow \mathbb{R}_+$

$$D_R(W \| W') = R(W) - R(W') - \langle \nabla R(W'), W - W' \rangle$$

- in Algo 1, initialize  $W_1 = \frac{1}{d} I_d$

$$W_{t+1} = \arg \min_{W \in W} \{ \eta \langle W, \tilde{L}_t \rangle + D_R(W \| W^t) \}$$

$$W \subset S \text{ with } \text{tr}(W) = 1$$

- choose  $R(W) = -\log \det(W)$

another aside

$$\frac{\partial (-\log \det W)}{\partial W_{ij}} = -\frac{1}{\det(W)} \cdot \frac{\partial \det W}{\partial W_{ij}} = -\frac{1}{\det(W)} \cdot \text{adj}(W)_{ij} = -(W^{-1})_{ij}$$

$$\text{then, } D_R(W \| U) = -\log \det(W) + \log \det(U) + \underbrace{\langle U^{-1}, W - U \rangle}_{= \text{tr}(U^{-1}W) - \text{tr}(I)}$$

$$= -\log \det(W) + \log \det U + \text{tr}(U^{-1}W) - d$$

$$= \log \det(UW^{-1}) + \text{tr}(U^{-1}W) - d$$

$$= \text{tr}(U^{-1}W) - \log \det(U^{-1}W) - d \quad \square$$

- Algo 1 iterates can be rewritten as

$$\tilde{W}_{t+1} = \arg \min_W \{ \eta \text{tr}(W \tilde{L}_t) + D_R(W \| W_t) \} \quad (\text{update})$$

$$W_{t+1} = \arg \min_{W \in \mathcal{W}} D_R(W \| \tilde{W}_{t+1})$$

-  $\tilde{W}_{t+1}$  satisfies (not verified yet)

$$\nabla R(\tilde{W}_{t+1}) = \nabla R(W_t) - \eta \tilde{L}_t$$

$$\Rightarrow \tilde{W}_{t+1}^{-1} = W_t^{-1} + \eta \tilde{L}_t$$

$$\Rightarrow \tilde{W}_{t+1} = (W_t^{-1} + \eta \tilde{L}_t)^{-1} = W_t^{1/2} \left( I + \eta W_t^{1/2} \tilde{L}_t W_t^{1/2} \right)^{-1} W_t^{1/2}$$

Lemma 9: (from an older reference) for any  $\eta > 0$ ,  $\gamma \in [0, 1]$

$$R_T \leq \frac{d \log T}{\eta} + 2\gamma T + 2 + (1-\gamma) \underbrace{\sum_{t=1}^T E[\langle W_t - \tilde{W}_{t+1}, \tilde{L}_t \rangle]}_{T\epsilon}$$

let  $B_t = W_t^{1/2} \tilde{L}_t W_t^{1/2}$ , then

$$\tilde{W}_{t+1} = W_t^{1/2} (I + \eta B_t)^{-1} W_t^{1/2} = W_t - \eta W_t^{1/2} \underbrace{B_t (I + \eta B_t)^{-1}}_{\tilde{L}_t} W_t^{1/2}$$

$$\begin{aligned} \text{use } (I+A)^{-1} &= I + (I+A)^{-1} - I \\ &= I + (I+A)^{-1} - (I+A)^{-1} (I+A) \\ &= I + (I - I - A)(I+A)^{-1} \\ &= I - A(I+A)^{-1} \end{aligned}$$

$$\langle W_t - W_{t+1}, L_t \rangle = \eta \text{tr} \left( W_t^{1/2} B_t (I + \eta B_t)^{-1} W_t^{1/2} \tilde{L}_t \right)$$

$$- \eta \operatorname{tr} (B_t (I + \eta B_t)^{-1} B_t)$$

$$= \eta \sum_{t=1}^d \frac{b_{ti}^2}{1 + \eta b_{ti}}$$

$b_{ti}$  - eig vals of  $B_t$

- these are bounded separately for dense, sparse cases
- lemma 10 bounds  $\operatorname{Tr} L_t \leq \eta (d^2 + 1) L_t^2$  (dense)
- lemma 11 bounds  $\operatorname{Tr} L_t \leq 8nd \|L_t\|_F^2$  (sparse)

Running time: - Dense -  $\tilde{O}(d^3)$  per trial  
 - sparse -  $\tilde{O}(d)$  per trial

### Possible extensions

- what about non-square loss matrices?
- does looking at  $w_t \in \mathbb{R}^{d \times r}$  make sense?
- can we look at sparse PCA?
- some other structure on  $w_t / L_t \dots$ ?
- missing data?? ie if we only see  $\operatorname{tr}(w_t w_t^T P_\Omega(L_t))$ ?  
 or other places for  $P_\Omega(\cdot)$  as well
- check contaminated data paper by Nowak grp for "good" models.