

1. Supplementary material

Lemma 1 (Lemma 14 in Appendix) *Let G_1 and G_2 be two zero mean Gaussian distributions with covariance matrix $\mathbf{\Gamma}\mathbf{\Sigma}\mathbf{\Gamma}$ and $\mathbf{\Gamma}\mathbf{\Theta}\mathbf{\Gamma}$. Furthermore $\mathbf{\Sigma}$ and $\mathbf{\Theta}$ are positive definite matrices. If there exists (i, j) such that*

$$|\mathbf{\Sigma}_{i,j} - \mathbf{\Theta}_{i,j}| \geq \delta(\mathbf{\Sigma}_{i,i} + \mathbf{\Theta}_{i,i} + \mathbf{\Sigma}_{j,j} + \mathbf{\Theta}_{j,j}), \quad (1)$$

then the total variation distance between G_1 and G_2 is at least $\frac{1}{12e^{1/4}}\delta$.

Proof Given $\phi_1(u)$ and $\phi_2(u)$ as characteristic function of G_1 and G_2 respectively. Due to Lemma 2 in (Moridomi et al., 2018), we have

$$\int_x |G_1(x) - G_2(x)| dx \geq \max_{u \in \mathbb{R}^N} |\phi_1(u) - \phi_2(u)|, \quad (2)$$

so we only need to show the lower bound of $\max_{u \in \mathbb{R}^N} |\phi_1(u) - \phi_2(u)|$.

Then we set that characteristic function of G_1 and G_2 are $\phi_1(u) = e^{\frac{-1}{2}u^T \mathbf{\Gamma}^T \mathbf{\Sigma} \mathbf{\Gamma} u}$ and $\phi_2(u) = e^{\frac{-1}{2}u^T \mathbf{\Gamma}^T \mathbf{\Theta} \mathbf{\Gamma} u}$ respectively. Set that $\alpha_1 = (\mathbf{\Gamma}v)^T \mathbf{\Sigma} (\mathbf{\Gamma}v)$, $\alpha_2 = (\mathbf{\Gamma}v)^T \mathbf{\Theta} (\mathbf{\Gamma}v)$ and $\mathbf{\Gamma}u = \frac{\mathbf{\Gamma}v}{\sqrt{\alpha_1 + \alpha_2}}$. Moreover we denote that $\bar{v} = \mathbf{\Gamma}v$, for any $\bar{v} \in \mathbb{R}^V$, there exists $v \in \mathbb{R}^V$. $\bar{u} = \mathbf{\Gamma}u$ in the same way.

We need only give the lower bound of $\max_{u \in \mathbb{R}^N} |\phi_1(u) - \phi_2(u)|$.

Next we have that

$$\begin{aligned} & \max_{u \in \mathbb{R}^N} |\phi_1(u) - \phi_2(u)| \\ &= \max_{u \in \mathbb{R}^N} \left| e^{\frac{-1}{2}u^T \mathbf{\Gamma}^T \mathbf{\Sigma} \mathbf{\Gamma} u} - e^{\frac{-1}{2}u^T \mathbf{\Gamma}^T \mathbf{\Theta} \mathbf{\Gamma} u} \right| \\ &= \max_{u \in \mathbb{R}^V} \left| e^{\frac{-1}{2}(\mathbf{\Gamma}u)^T \mathbf{\Sigma} (\mathbf{\Gamma}u)} - e^{\frac{-1}{2}(\mathbf{\Gamma}u)^T \mathbf{\Theta} (\mathbf{\Gamma}u)} \right| \\ &\geq \max_{\bar{v} \in \mathbb{R}^N} \left| e^{\frac{-\alpha_1}{2(\alpha_1 + \alpha_2)}} - e^{\frac{-\alpha_2}{2(\alpha_1 + \alpha_2)}} \right| \\ &\geq \max_{\bar{v} \in \mathbb{R}^N} \left| \frac{1}{2e^{1/4}} \frac{\alpha_1 - \alpha_2}{\alpha_1 + \alpha_2} \right|. \end{aligned} \quad (3)$$

Then second inequality is due to Lemma 5, since $\min\{\frac{\alpha_1}{\alpha_1 + \alpha_2}, \frac{\alpha_2}{\alpha_1 + \alpha_2}\} \in (0, \frac{1}{2}]$.

Due to assumption in the Lemma we obtain for some (i, j) that

$$\begin{aligned} & \delta(\mathbf{\Sigma}_{i,i} + \mathbf{\Theta}_{i,i} + \mathbf{\Sigma}_{j,j} + \mathbf{\Theta}_{j,j}) \leq |\mathbf{\Sigma}_{i,j} - \mathbf{\Theta}_{i,j}| \\ &= \frac{1}{2} |(\mathbf{e}_i + \mathbf{e}_j)^T (\mathbf{\Sigma} - \mathbf{\Theta})(\mathbf{e}_i + \mathbf{e}_j) - \mathbf{e}_i^T (\mathbf{\Sigma} - \mathbf{\Theta}) \mathbf{e}_i - \mathbf{e}_j^T (\mathbf{\Sigma} - \mathbf{\Theta}) \mathbf{e}_j| \end{aligned} \quad (4)$$

It implies that one of $(\mathbf{e}_i + \mathbf{e}_j)^T (\mathbf{\Sigma} - \mathbf{\Theta})(\mathbf{e}_i + \mathbf{e}_j)$, $\mathbf{e}_i^T (\mathbf{\Sigma} - \mathbf{\Theta}) \mathbf{e}_i$ and $\mathbf{e}_j^T (\mathbf{\Sigma} - \mathbf{\Theta}) \mathbf{e}_j$ has absolute value greater than $\frac{2\delta}{3}(\mathbf{\Sigma}_{i,i} + \mathbf{\Theta}_{i,i} + \mathbf{\Sigma}_{j,j} + \mathbf{\Theta}_{j,j})$.

Since $\mathbf{\Sigma}, \mathbf{\Theta}$ are strictly positive definite matrices, we have that for all $v \in \{\mathbf{e}_i + \mathbf{e}_j, \mathbf{e}_i, \mathbf{e}_j\}$

$$v^T (\mathbf{\Sigma} + \mathbf{\Theta}) v \leq 2(\mathbf{\Sigma} + \mathbf{\Theta})_{i,i} + 2(\mathbf{\Sigma} + \mathbf{\Theta})_{j,j}. \quad (5)$$

and therefore we have that

$$\max_{\bar{v} \in \mathbb{R}^N} \left| \frac{1}{2e^{1/4}} \frac{\alpha_1 - \alpha_2}{\alpha_1 + \alpha_2} \right| \geq \max_{\bar{v} \in \{\mathbf{e}_i + \mathbf{e}_j, \mathbf{e}_i, \mathbf{e}_j\}} \left| \frac{1}{2e^{1/4}} \frac{v^T(\boldsymbol{\Sigma} - \boldsymbol{\Theta})v}{v^T(\boldsymbol{\Sigma} + \boldsymbol{\Theta})v} \right| \geq \frac{\delta}{6e^{1/4}} \quad (6)$$

■

Now we give the proof of the Main Theorem as follows:

Proof [Proof of Theorem 2] Due to Lemma 4 (in main part) we obtain that

$$\text{Regret}_{\text{OSDP}}(T, \mathcal{K}, \mathcal{L}, \mathbf{W}^*) \leq \frac{H_0}{\eta} + \frac{\eta}{s} T. \quad (7)$$

Due to the main proposition in main part we know that $s = 1/(1152(\beta + \rho\epsilon)^2\sqrt{e}g^2)$.

Thus we need only to show $H_0 \leq \frac{\tau}{\epsilon}$. Given \mathbf{W}_0 and \mathbf{W}_1 is the minimizer and maximizer of R respectively, then we obtain that

$$\begin{aligned} \max_{\mathbf{W}, \mathbf{W}' \in \mathcal{K}} (R(\mathbf{W}) - R(\mathbf{W}')) &= R(\mathbf{W}_1) - R(\mathbf{W}_0) \\ &= -\ln \det(\boldsymbol{\Gamma}\mathbf{W}_1\boldsymbol{\Gamma} + \epsilon\mathbf{E}) + \ln \det(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma} + \epsilon\mathbf{E}) \\ &= \sum_{i=1}^N \ln \frac{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma}) + \epsilon}{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_1\boldsymbol{\Gamma}) + \epsilon} \\ &= \sum_{i=1}^N \ln \left(\frac{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma})}{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_1\boldsymbol{\Gamma}) + \epsilon} + \frac{\epsilon}{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_1\boldsymbol{\Gamma}) + \epsilon} \right) \\ &\leq \sum_{i=1}^N \ln \left(\frac{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma})}{\epsilon} + 1 \right) \\ &\leq \sum_{i=1}^N \frac{\lambda_i(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma})}{\epsilon} = \frac{\text{Tr}(\boldsymbol{\Gamma}\mathbf{W}_0\boldsymbol{\Gamma})}{\epsilon} \leq \frac{\tau}{\epsilon}. \end{aligned} \quad (8)$$

Plugging s , we obtain that

$$\text{Regret}_{\text{OSDP}}(T, \mathcal{K}, \mathcal{L}, \mathbf{W}^*) = O \left(g^2(\beta + \rho\epsilon)^2 T \eta + \frac{\tau}{\epsilon \eta} \right). \quad (9)$$

■

Lemma 2 (Lemma A.1 (Moridomi et al., 2018)) Let P and Q be probability distributions over \mathbb{R}^N and $\phi_P(u)$ and $\phi_Q(u)$ be their characteristic functions, respectively. Then

$$\max_{u \in \mathbb{R}^N} |\phi_P(u) - \phi_Q(u)| \leq \int_x |P(x) - Q(x)| dx, \quad (10)$$

the right hand side is the total variation distance between any distribution Q and P .

Lemma 3 (Lemma A.2 (Christiano, 2014)) *Let P and Q be probability distributions over \mathbb{R}^N with total variation distance δ . Then*

$$H(\alpha P + (1 - \alpha)Q) \leq \alpha H(P) + (1 - \alpha)H(Q) - \alpha(1 - \alpha)\delta^2, \quad (11)$$

where $H(P) = \mathbb{E}_{x \sim P}[\ln P(x)]$.

Lemma 4 (Lemma A.3 (Moridomi et al., 2018)) *For any probability distribution P over \mathbb{R}^N with zero mean and covariance matrix Σ , its entropy is bounded by the log-determinant of covariance matrix. That is*

$$-H(P) \leq \frac{1}{2} \ln(\det(\Sigma)(2\pi e)^N). \quad (12)$$

Lemma 5 (Lemma A.4 (Moridomi et al., 2018))

$$e^{-\frac{x}{2}} - e^{-\frac{1-x}{2}} \geq \frac{e^{-1/4}}{2}(1 - 2x), \quad (13)$$

for $0 \leq x \leq 1/2$.

2. Definition of biclustered structure and ideal quasi dimension

As in Herbster et al. (2020), we define the class of (k, l) -biclustered structure matrices as follows:

Definition 6 *For $m \geq k$ and $n \geq l$, the class of (k, l) -binary biclustered matrices is defined as*

$$\mathbb{B}_{k,l}^{m \times n} = \{\mathbf{U} \in \{-1, +1\}^{m \times n} : \mathbf{r} \in [k]^m, \mathbf{c} \in [l]^n, \mathbf{V} \in \{1, -1\}^{k \times l}, \mathbf{U}_{i,j} = \mathbf{V}_{r_i, c_j}, i \in [m], j \in [n]\}.$$

Denote $\mathcal{B}^{m,d} = \{\mathbf{R} \subset \{0, 1\}^{m \times d} : \|\mathbf{R}_i\|_2 = 1, i \in [m], \text{rank}(\mathbf{R}) = d\}$, for any matrix $\mathbf{U} \in \mathbb{B}_{k,l}^{m,n}$ we can decompose $\mathbf{U} = \mathbf{R}\mathbf{U}^*\mathbf{C}^\top$ for some $\mathbf{U}^* \in \{-1, +1\}^{k \times l}$, $\mathbf{R} \in \mathcal{B}^{m,k}$ and $\mathbf{C} \in \mathcal{B}^{n,l}$.

Theorem 7 ((Herbster et al., 2020)) *If $\mathbf{U} \in \mathbb{B}_{k,l}^{m \times n}$ define $\mathcal{D}_{\mathbf{M}, \mathbf{N}}^o(\mathbf{U})$ as*

$$\mathcal{D}_{\mathbf{M}, \mathbf{N}}^o(\mathbf{U}) = 2\text{Tr}(\mathbf{R}^\top \mathbf{M} \mathbf{R})\alpha_{\mathbf{M}} + 2\text{Tr}(\mathbf{C}^\top \mathbf{N} \mathbf{C})\alpha_{\mathbf{N}} + 2k + 2l, \quad (14)$$

where \mathbf{M}, \mathbf{N} are PD-Laplacian, as the minimum over all decompositions of $\mathbf{U} = \mathbf{R}\mathbf{U}^*\mathbf{C}^\top$ for some $\mathbf{U}^* \in \{-1, +1\}^{k \times l}$, $\mathbf{R} \in \mathcal{B}^{m,k}$ and $\mathbf{C} \in \mathcal{B}^{n,l}$. Thus, for $\mathbf{U} \in \mathbb{B}_{k,l}^{m \times n}$,

$$\mathcal{D}_{\mathbf{M}, \mathbf{N}}^\gamma(\mathbf{U}) \leq \mathcal{D}_{\mathbf{M}, \mathbf{N}}^o(\mathbf{U}), \quad (15)$$

if $\|\mathbf{U}\|_{\max} \leq \frac{1}{\gamma}$.

Moreover, we define the max-norm of a matrix $\mathbf{U} \in \mathbb{R}^{m \times n}$ as follows:

$$\|\mathbf{U}\|_{\max} = \min_{\mathbf{P}\mathbf{Q}^\top = \mathbf{U}} \left\{ \max_{1 \leq i \leq m} \|\mathbf{P}_i\| \max_{1 \leq j \leq n} \|\mathbf{Q}_j\| \right\}. \quad (16)$$

Furthermore we define the quasi-dimension of a matrix \mathbf{U} with $\mathbf{M} \in \mathbb{S}_{++}^{m \times m}$ and $\mathbf{N} \in \mathbb{S}_{++}^{n \times n}$ at margin γ as

$$\mathcal{D}_{\mathbf{M}, \mathbf{N}}^\gamma(\mathbf{U}) = \min_{\bar{\mathbf{P}}\bar{\mathbf{Q}}^\top = \gamma\mathbf{U}} \alpha_M \text{Tr}(\bar{\mathbf{P}}^\top \mathbf{M} \bar{\mathbf{Q}}) + \alpha_N \text{Tr}(\bar{\mathbf{Q}}^\top \mathbf{N} \bar{\mathbf{Q}}). \quad (17)$$

See section 4.1 from [Herbster et al. \(2020\)](#), if \mathbf{U} is a (k, l) -biclustered structured matrix, they show an example where $\mathcal{D}_{\mathbf{M}, \mathbf{N}}^\gamma(\mathbf{U}) \in O(k + l)$ with ideal side information. When exactly that there exists a sequence that $y_t = (\bar{\mathbf{P}}\bar{\mathbf{Q}}^\top)_{i_t, j_t} = \mathbf{U}_{i_t, j_t}$ where $(\bar{\mathbf{P}}, \bar{\mathbf{Q}}) = \arg \min_{\mathbf{P}, \mathbf{Q}} \mathcal{D}_{\mathbf{M}, \mathbf{N}}^\gamma(\mathbf{U})$, and \mathbf{U} satisfies the assumptions in [Herbster et al. \(2020\)](#), then we have that $\hat{\mathbf{D}} \in O(k + l)$ with same side information.

3. Online similarity prediction with side information

In this section, we show the application of our reduction method and generalised log-determinant regularizer to online similarity prediction with side information.

Let $G = (V, E)$ be an undirected and connected graph with $n = |V|$ vertices and $m = |E|$ edges. Assign vertices to K classes with a vector $\mathbf{y} = \{y_1, \dots, y_n\}$ where $y_i \in \{1, \dots, K\}$. For a matrix \mathbf{L} , we denote \mathbf{L}^+ as pseudo-inverse matrix of \mathbf{L} . The online similarity prediction is defined as follows: On each round t , for a given pair of vertices (i_t, j_t) algorithm needs to predict whether they are in the same class denoted as \hat{y}_{i_t, j_t} . If they are in the same class then $y_{i_t, j_t} = 1$, $y_{i_t, j_t} = -1$, otherwise. Our target is to give a bound of the prediction mistakes $M = \sum_{t=1}^T \mathbb{1}_{\hat{y}_{i_t, j_t} \neq y_{i_t, j_t}}$.

Definition 8 *The set of cut-edges in (G, \mathbf{y}) is denoted as $\Phi^G(\mathbf{y}) = \{(i, j) \in E : y_i \neq y_j\}$ we abbreviate it to Φ^G and the cut-size is given as $|\Phi^G(\mathbf{y})|$. The set of cut-edges with respect to class label k is denoted as $\Phi_k^G(\mathbf{y}) = \{(i, j) \in E : k \in \{y_i, y_j\}, y_i \neq y_j\}$. Note that $\sum_{s=1}^k |\Phi_s^G(\mathbf{y})| = 2|\Phi^G(\mathbf{y})|$. Given $\mathbf{A} \in \mathbb{R}^{n \times n}$ such that $\mathbf{A}_{ij} = \mathbf{A}_{ji} = 1$ if $(i, j) \in E(G)$ and $\mathbf{A}_{ij} = 0$, otherwise. \mathbf{D} is denoted as diagonal matrix with \mathbf{D}_{ii} is the degree of vertex i . We define the Laplacian as $\mathbf{L} = \mathbf{D} - \mathbf{A}$.*

Definition 9 *If G is identified with a resistive network such that each edge is a unit resistor, then the effective resistance $R_{i,j}^G$ between pair $(i, j) \in V^2$ can be defined as $R_{i,j}^G = (e_i - e_j)\mathbf{L}^+(e_i - e_j)$, where e_i is the i -th vector in the canonical basis of \mathbb{R}^n .*

[Gentile et al. \(2013\)](#) gave a mistake bound in the following proposition:

Proposition 10 *Let (G, \mathbf{y}) be a labeled graph. If we run the Matrix Winnow with G as input graph, we have the following mistake bound*

$$M^W = O \left(|\Phi^G| \max_{(i,j) \in V^2} R_{i,j}^G \ln n \right) \quad (18)$$

In our new reduction, we define the comparator matrix $\mathbf{U} \in \{1, -1\}^{n \times n}$ where if vertices i, j are in the same class then $\mathbf{U}_{ij} = 1$, and $\mathbf{U}_{ij} = -1$, otherwise. Firstly, we denote that $\mathbf{1}$ is a K -dimensional vector that all entries are 1. Due to (Gentile et al., 2013; Herbster et al., 2020), we see that \mathbf{U} is a (K, K) -biclustered $n \times n$ matrix where $\mathbf{U}^* = 2\mathbf{I}_K - \mathbf{1}\mathbf{1}^\top$, and there exists $\mathbf{R} \in \mathcal{B}^{n, k}$ such that $\mathbf{U} = \mathbf{R}\mathbf{U}^*\mathbf{R}^\top$. Define the side information matrices $\mathbf{M} = \mathbf{N} \in \mathbb{R}^{n \times n}$ as PD-Laplacian $\tilde{\mathbf{L}}$, where \mathbf{L} is the Laplacian matrix based on the graph G .

Thus we have

$$\mathbf{\Gamma} = \begin{bmatrix} \sqrt{\alpha_{\tilde{\mathbf{L}}}} & 0 \\ 0 & \sqrt{\alpha_{\tilde{\mathbf{L}}}} \end{bmatrix}, \quad (19)$$

where $\alpha_{\tilde{\mathbf{L}}} = \max_i(\tilde{\mathbf{L}}^{-1})_{ii}$.

According to Herbster et al. (2020), we further obtain that $\frac{1}{\gamma} \in O(1)$, more concisely we can set that $\frac{1}{\gamma} = 3$. Meanwhile given sparse matrix \mathbf{Z} in the following equation

$$\mathbf{Z}\langle i, j \rangle = \frac{1}{2}(\mathbf{e}_i \mathbf{e}_{n+j}^\top + \mathbf{e}_{n+j} \mathbf{e}_i^\top). \quad (20)$$

Thus we give the following proposition for our reduction from a graph based online similarity prediction to a generalised OSDP problem $(\mathcal{K}, \mathcal{L})$ with bounded $\mathbf{\Gamma}$ -trace norm.

Proposition 11 *Given an online similarity prediction problem with graph (G, \mathbf{y}) , then we can reduce this problem to a generalised OSDP problem $(\mathcal{K}, \mathcal{L})$ with bounded $\mathbf{\Gamma}$ -trace norm such that*

$$\begin{aligned} \mathcal{K} &= \left\{ \mathbf{X} \in \mathbb{S}_{++}^{n \times n} : |\mathbf{X}_{ii}| \leq 1, \text{Tr}(\mathbf{\Gamma}\mathbf{X}\mathbf{\Gamma}) \leq \widehat{\mathcal{D}} \right\} \\ \mathcal{L} &= \{c\mathbf{Z}\langle i, j \rangle : c \in \{-1/\gamma, 1/\gamma\}, i \in [n], j \in [n]\}, \end{aligned}$$

where $\mathbf{\Gamma}$ is defined as above, and $\widehat{\mathcal{D}}$ is arbitrary. In particular, we have that

$$M = \sum_{t=1}^T \mathbb{I}_{\hat{y}_{i_t, j_t} \neq y_{i_t, j_t}} \leq \text{Regret}_{\text{OSDP}}(M, \mathcal{K}, \mathcal{L}) \quad (21)$$

According to Herbster et al. (2020), there exists $\bar{\mathbf{P}}, \bar{\mathbf{Q}}$ such that $\mathbf{U}^* = \bar{\mathbf{P}}\bar{\mathbf{Q}}^\top$, it implies that the hinge loss $\text{hloss}(\mathcal{S}, \gamma) = 0$.

Remark 12 *According to Theorem 3 and section 4.2 in (Herbster et al., 2020) if \mathbf{U} obtains the (K, K) -biclustered structure, i.e., there exists \mathbf{U}^* , such that $\mathbf{U}^* = 2\mathbf{I}_K - \mathbf{1}\mathbf{1}^\top$, and there exists $\mathbf{R} \in \mathcal{B}^{n, k}$ such that $\mathbf{U} = \mathbf{R}\mathbf{U}^*\mathbf{R}^\top$, due to our Theorem 13 in main part, we have that*

$$M \leq O\left(\text{Tr}(\mathbf{R}^\top \mathbf{L} \mathbf{R}) \alpha_{\mathbf{L}}\right), \quad (22)$$

where \mathbf{L} is Laplacian of the corresponding graph G .

Remark 13 According to [Herbster et al. \(2020\)](#), we have that

$$\text{Tr}(\mathbf{R}^\top \mathbf{L} \mathbf{R}) \leq 2 \sum_{(i,j) \in E} \|\mathbf{R}_i - \mathbf{R}_j\|^2,$$

where $\sum_{(i,j) \in E} \|\mathbf{R}_i - \mathbf{R}_j\|^2$ counts only when there is an edge between different classes. Due to the definition of $|\Phi^G|$, we have that $\sum_{(i,j) \in E} \|\mathbf{R}_i - \mathbf{R}_j\|^2 = |\Phi^G|$.

Simultaneously, $\alpha_{\mathbf{L}} = \max_{i \in [n]} \mathbf{L}_{ii}^+$ so we obtain that $\alpha_{\mathbf{L}} \geq \mathbf{e}_i^\top \mathbf{L}^+ \mathbf{e}_i, \forall i \in [n]$. It implies that $4\alpha_{\mathbf{L}} \geq \max_{(i,j) \in V^2} R_{i,j}^G$. Thus we have that our new mistake bound improves the previous bound a logarithmic factor and recovers the previous bound up to a constant factor.

References

- Paul Christiano. Online local learning via semidefinite programming. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 468–474. ACM, 2014.
- Claudio Gentile, Mark Herbster, and Stephen Pasteris. Online similarity prediction of networked data from known and unknown graphs. In *Conference on Learning Theory*, pages 662–695, 2013.
- Mark Herbster, Stephen Pasteris, and Lisa Tse. Online matrix completion with side information. *Advances in Neural Information Processing Systems*, 33, 2020.
- Ken-ichiro Moridomi, Kohei Hatano, and Eiji Takimoto. Online linear optimization with the log-determinant regularizer. *IEICE Transactions on Information and Systems*, 101(6):1511–1520, 2018.