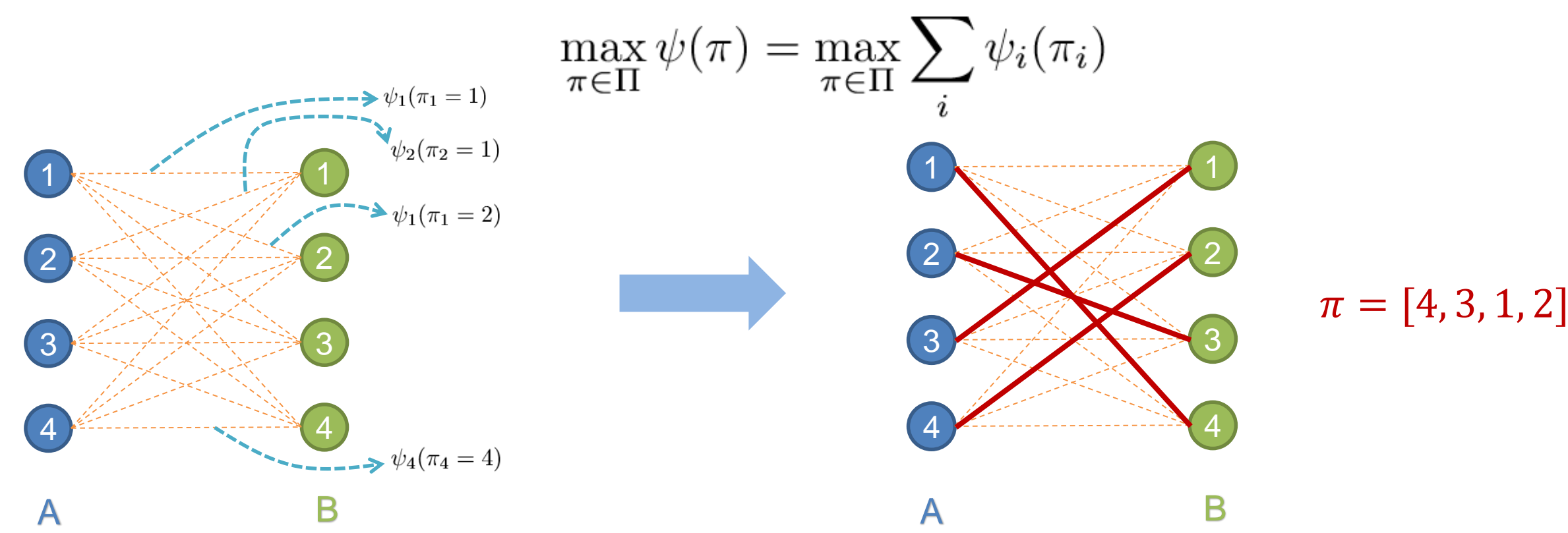


Bipartite Matching Task

Weighted Bipartite Matching

- Given: ① two sets of elements A and B with equal size, ② weights between the elements in A and the elements in B
- Task: find one-to-one mapping that maximize sum of potentials:



Learning Bipartite Matching Task

- Given: training data \rightarrow each sample: a bipartite graph (x) and a ground truth assignment (π)
- Task: learn weight function $\psi_i(\cdot)$ that minimizes miss-assignment metric (e.g. Hamming loss)

Applications:

- Word alignment (Taskar et. al., 2005; Pado & Lapta, 2006; Mac-Cartney et. al., 2008)
- Correspondence between images (Belongie et. al., 2002; Dellaert et. al., 2003)
- Learning to rank documents (Dwork et. al., 2001; Le & Smola, 2007)

Previous Methods and Shortcomings

Desiderata for a Predictor:

- Efficiency: learning & prediction runtime is in a (low degree) polynomial time
- Consistency: must also minimize Hamming loss under ideal condition (given the true distribution and fully expressive model parameters)

① Exponential Family Random Field Approach (Petterson et. al., 2009; Volkovs & Zemel, 2012)

Probabilistic model:

$$P_\psi(\pi) = \frac{1}{Z_\psi} \exp\left(\sum_{i=1}^n \psi_i(\pi_i)\right)$$

$$Z_\psi = \sum_{\pi} \prod_{i=1}^n \exp(\psi_i(\pi_i)) = \text{perm}(\mathbf{M}) \quad \text{where } M_{i,j} = \exp(\psi_i(j))$$



Consistent? Yes!

produce Bayes optimal prediction over the Hamming loss in an ideal condition



Efficient? No!

normalization term Z_ψ involves matrix permanent computation (a #P-hard problem)

② Maximum Margin Approach (Tschantaridis et. al., 2005)

Max-margin model:

$$\min_{\psi} \mathbb{E}_{\pi \sim \tilde{P}} \left[\max_{\pi'} \{ \text{loss}(\pi, \pi') + \psi(\pi') \} - \psi(\pi) \right]$$

\tilde{P} is the empirical distribution



Efficient? Yes!

polynomial algorithm for computing the maximum violated constraint (Hungarian algorithm)



Consistent? No!

based on the CS multiclass SVM: not consistent for distributions with no majority label

Adversarial Bipartite Matching

Formulation

- Our method seeks a predictor that robustly minimizes Hamming loss, against the worst-case permutation mixture probability that is consistent with the statistics of the training data

$$\min_{\tilde{P}} \max_{\hat{P}} \mathbb{E}_{x \sim \tilde{P}; \tilde{\pi} \sim \hat{P}; \hat{\pi} \sim \tilde{P}} [\text{loss}(\hat{\pi}, \tilde{\pi})]$$

$$\text{s.t. } \mathbb{E}_{x \sim \tilde{P}; \tilde{\pi} \sim \hat{P}} \left[\sum_{i=1}^n \phi_i(x, \tilde{\pi}_i) \right] = \mathbb{E}_{(x, \pi) \sim \tilde{P}} \left[\sum_{i=1}^n \phi_i(x, \pi_i) \right]$$

- Predictor: - makes a probabilistic prediction $\hat{P}(\hat{\pi}|x)$ and aims to minimize the loss - is pitted with an adversary instead of the empirical distribution
- Adversary: - makes a probabilistic prediction $\tilde{P}(\tilde{\pi}|x)$ and aims to maximize the loss - constrained to select probability that match the statistics of empirical distribution (\tilde{P}) via moment matching on the features $\phi(x, \pi) = \sum_{i=1}^n \phi_i(x, \pi_i)$

Dual Formulation

$$\min_{\theta} \mathbb{E}_{x, \pi \sim \tilde{P}} \min_{\tilde{P}} \max_{\hat{P}} \mathbb{E}_{\tilde{\pi} \sim \hat{P}; \hat{\pi} \sim \tilde{P}} \left[\text{loss}(\hat{\pi}, \tilde{\pi}) + \theta \cdot \sum_{i=1}^n (\phi_i(x, \tilde{\pi}_i) - \phi_i(x, \pi_i)) \right]$$

where θ is the dual variable for moment matching constraints

Augmented Hamming loss matrix for $n = 3$ permutations:

	$\tilde{\pi} = 123$	$\tilde{\pi} = 132$	$\tilde{\pi} = 213$	$\tilde{\pi} = 231$	$\tilde{\pi} = 312$	$\tilde{\pi} = 321$
$\hat{\pi} = 123$	0 + δ_{123}	2 + δ_{132}	2 + δ_{213}	3 + δ_{231}	3 + δ_{312}	2 + δ_{321}
$\hat{\pi} = 132$	2 + δ_{123}	0 + δ_{132}	3 + δ_{213}	2 + δ_{231}	2 + δ_{312}	3 + δ_{321}
$\hat{\pi} = 213$	2 + δ_{123}	3 + δ_{132}	0 + δ_{213}	2 + δ_{231}	2 + δ_{312}	3 + δ_{321}
$\hat{\pi} = 231$	3 + δ_{123}	2 + δ_{132}	2 + δ_{213}	0 + δ_{231}	3 + δ_{312}	2 + δ_{321}
$\hat{\pi} = 312$	3 + δ_{123}	2 + δ_{132}	2 + δ_{213}	3 + δ_{231}	0 + δ_{312}	2 + δ_{321}
$\hat{\pi} = 321$	2 + δ_{123}	3 + δ_{132}	3 + δ_{213}	2 + δ_{231}	2 + δ_{312}	0 + δ_{321}

size: $n! \times n!$

Intractable for modestly-sized n

Algorithms

① Double Oracle Method

- Based on the observation: equilibrium is usually supported by small number of permutations
- Iterative method: - start from a single permutation for each player - alternately: * compute predictor's (/adversary's) strategy in the current game * compute adversary's (/predictor's) best response, add to the game - until no improvement in the game value
- Use the game solution to compute the gradient and perform gradient update
- No formal polynomial bound is known \rightarrow the whole runtime cannot be characterized as polynomial

② Marginal Distribution Formulation

- Reformulation: from permutation mixture distributions to marginal distributions:

$$\mathbf{P} = \begin{matrix} & \text{Predictor} \\ & \begin{matrix} 1 & 2 & 3 \\ \tilde{\pi}_1 & p_{1,1} & p_{1,2} & p_{1,3} \\ \tilde{\pi}_2 & p_{2,1} & p_{2,2} & p_{2,3} \\ \tilde{\pi}_3 & p_{3,1} & p_{3,2} & p_{3,3} \end{matrix} \\ \mathbf{P} = & \end{matrix}$$

$$\mathbf{Q} = \begin{matrix} & \text{Adversary} \\ & \begin{matrix} 1 & 2 & 3 \\ \hat{\pi}_1 & q_{1,1} & q_{1,2} & q_{1,3} \\ \hat{\pi}_2 & q_{2,1} & q_{2,2} & q_{2,3} \\ \hat{\pi}_3 & q_{3,1} & q_{3,2} & q_{3,3} \end{matrix} \\ \mathbf{Q} = & \end{matrix}$$

$$p_{i,j} = \hat{P}(\hat{\pi}_i = j) \quad q_{i,j} = \tilde{P}(\tilde{\pi}_i = j)$$

Birkhoff – Von Neumann theorem:

The convex hull of the set of permutations forms a convex polytope whose points are doubly stochastic matrices: $\mathbf{P}\mathbf{1} = \mathbf{P}^T\mathbf{1} = \mathbf{Q}\mathbf{1} = \mathbf{Q}^T\mathbf{1} = \mathbf{1}$

- Reduce the space of optimization from $O(n!)$ to $O(n^2)$
- Marginal optimization (after adding regularization and smoothing penalty):

$$\max_{\mathbf{Q} \geq 0} \min_{\theta} \frac{1}{m} \sum_{i=1}^m \min_{\mathbf{P}_i \geq 0} \left[\langle \mathbf{Q}_i - \mathbf{Y}_i, \sum_k \theta_k \mathbf{X}_{i,k} \rangle - \langle \mathbf{P}_i, \mathbf{Q}_i \rangle + \frac{\mu}{2} \|\mathbf{P}_i\|_F^2 - \frac{\mu}{2} \|\mathbf{Q}_i\|_F^2 \right] + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\text{s.t. } \mathbf{P}_i \mathbf{1} = \mathbf{P}_i^T \mathbf{1} = \mathbf{Q}_i \mathbf{1} = \mathbf{Q}_i^T \mathbf{1} = \mathbf{1}, \quad \forall i$$

- Techniques: - Outer (Q): projected Quasi-Newton with a projection to doubly-stochastic matrix - Inner (θ): closed-form solution - Inner (P): projection to doubly-stochastic matrix - Projection to doubly-stochastic matrix: ADMM

Consistency

Empirical Risk Perspective of Adversarial Bipartite Matching

- Adversarial Bipartite Matching can be viewed as an ERM method with surrogate loss AL_f^{perm}

$$\min_{\theta} \mathbb{E}_{x \sim \tilde{P}} \max_{\pi \sim \tilde{P}} [AL_{f_\theta}^{\text{perm}}(x, \pi)]$$

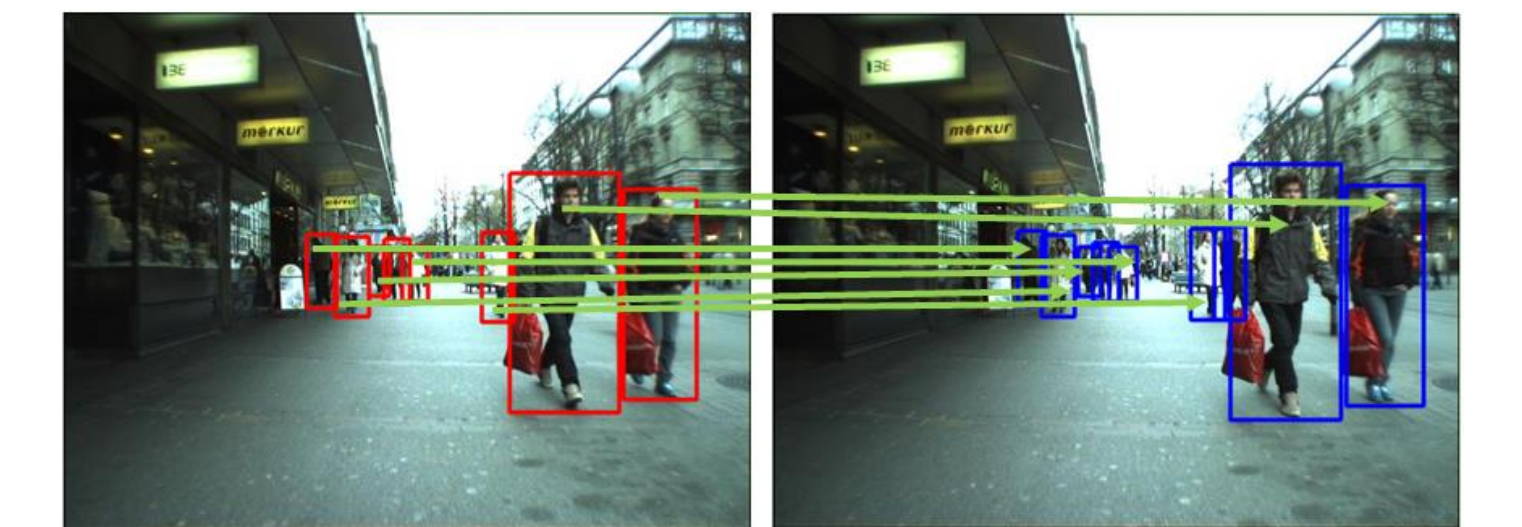
$$\text{where: } AL_{f_\theta}^{\text{perm}}(x, \pi) \triangleq \min_{\tilde{P}(\tilde{\pi}|x)} \max_{\hat{P}(\hat{\pi}|x)} \mathbb{E}_{\tilde{\pi} \sim \tilde{P}; \hat{\pi} \sim \hat{P}} [\text{loss}(\hat{\pi}, \tilde{\pi}) + f_\theta(x, \tilde{\pi}) - f_\theta(x, \pi)]$$

- We show that minimizing AL_f^{perm} also minimizes the Hamming loss given true distribution, and f is optimized over the set of all measurable functions on the input space (x, π)
- The consistency result also holds when f is an additive function over individual assignment π_i

Experiments

Experiments | Video Tracking Tasks

- Predict object correspondence in two different frames
- Public benchmark datasets
- 5 datasets in 2 groups (TUD and ETH)
- 48 different features for each pair of objects
- Train on one dataset, test on another dataset from the same group



Experiment results

Table 4. The mean and standard deviation (in parenthesis) of the average accuracy (1 - the average Hamming loss) for the adversarial bipartite matching model compared with Structured-SVM.

TRAINING/TESTING	ADV DO	ADV MARG.	SSVM	ADV DO #PERM.
CAMPUS/STADTMITTE	0.662 (0.09)	0.662 (0.08)	0.662 (0.08)	11.4
STADTMITTE/CAMPUS	0.672 (0.12)	0.667 (0.11)	0.660 (0.12)	7.4
BAHNHOF/SUNNYDAY	0.758 (0.12)	0.754 (0.10)	0.729 (0.15)	5.8
PEDCROSS2/SUNNYDAY	0.760 (0.08)	0.750 (0.10)	0.736 (0.13)	8.2
SUNNYDAY/BAHNHOF	0.755 (0.20)	0.751 (0.18)	0.739 (0.20)	9.8
PEDCROSS2/BAHNHOF	0.760 (0.12)	0.763 (0.16)	0.731 (0.21)	10.8
BAHNHOF/PEDCROSS2	0.718 (0.16)	0.714 (0.16)	0.701 (0.18)	8.5
SUNNYDAY/PEDCROSS2	0.719 (0.18)	0.712 (0.17)	0.700 (0.18)	14.4

Empirical runtime (until convergence)

Table 5. Running time (in seconds) of the model for various number of elements n with fixed number of samples ($m = 50$)

DATASET	# ELEMENTS	ADV MARG.	SSVM
CAMPUS	12	1.96	0.22
STADTMITTE	16	2.46	0.25
SUNNYDAY	18	2.75	0.15
PEDCROSS2	30	8.18	0.26
BAHNHOF	34	9.79	0.31

Adv. Marginal Formulation: grows (roughly) quadratically in n

Adv. Double Oracle: small number of permutations in the equilibrium

Adversarial Bipartite Matching:

6 pairs of datasets: significantly outperforms SSVM

2 pairs of datasets: competitive with SSVM

Conclusions

Exponential Family Random Field

(Petterson et. al., 2009; Volkovs & Zemel, 2012)

Efficient? Consistent? Perform well?

✗ ✓ ?

Maximum Margin

(Tschantaridis et. al., 2005)

✓ ✗ ??

Adversarial Bipartite Matching

(our approach)

✓ ✓ ✓

Acknowledgement

This research was supported in part by NSF Grants RI-1526379 and CAREER-#1652530.