

Adversarial Surrogate Losses for General Multiclass Classification

Rizal Zaini Ahmad Fathony

Committee: Prof. Brian Ziebart (Chair) Prof. Bhaskar DasGupta Prof. Lev Reyzin Prof. Xinhua Zhang Prof. Simon Lacoste-Julien

Supervised Learning \rightarrow Multiclass Classification



Multiclass Classification \rightarrow Zero-One Loss

Digit Recognition



Loss Function: $loss(\hat{y}, y) = I(\hat{y} \neq y)$

General Multiclass Classification \rightarrow any loss

Multiclass Classification \rightarrow Ordinal Classification

Movie Rating Prediction









Predicted vs Actual Label:



Loss Function (example): $loss(\hat{y}, y) = |\hat{y} - y|$

Multiclass Classification \rightarrow Taxonomy Classification

Object Classification



Loss Function (example): $loss(\hat{y}, y) = h - v(\hat{y}, y) + 1$

h : tree height $v(\hat{y}, y)$: level of the common ancestor loss(Cat,Dog) = 1 loss(Cat,Cow) = 2 loss(Cow,Person) = 3 Loss(Cow,Motorbike)= 4

loss(Bus,Car) = 2
loss(Bus,Bicycle) = 3
loss(Car,Cow) = 4
loss(Bus,Person) = 4

Multiclass Classification \rightarrow Loss Matrix

loss function: $loss(\hat{y}, y) \rightarrow loss$ matrix: L

Zero One Loss $loss(\hat{y}, y) = I(\hat{y} \neq y)$

 $\begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix}$

Zero One Loss Ordinal Classification Loss

$$loss(\hat{y}, y) = |\hat{y} - y|$$

0	1	2	3	4
1	0	1	2	3
2	1	0	1	2
3	2	1	0	1
4	3	2	1	0

Taxonomy-based loss

 $loss(\hat{y}, y) = h - v(\hat{y}, y) + 1$



Empirical Risk Minimization (ERM)

- Assume a family of parametric hypothesis function *f* (e.g. linear discriminator)
- Find the hypothesis f^* that minimize the empirical risk:

$$\min_{f} \frac{1}{n} \sum_{i=1}^{n} \operatorname{loss}(f(\mathbf{x}_{i}), y_{i}) = \min_{f} \mathbb{E}_{\tilde{P}(\mathbf{x}, y)} \left[\operatorname{loss}(f(\mathbf{X}), Y) \right]$$

Intractable optimization, non-convex, non-continuous

Convex surrogate loss need to be employed

Example:

Binary zero-one loss

Surrogate Loss:

- Hinge loss (used by SVM)
- Log loss (used by Logistic Regression)
- Exponential loss (used by AdaBoost)



ERM under Hinge Loss and Log Loss

SVM (hinge loss):

 $\min_{f} \mathbb{E}_{\tilde{P}(\mathbf{x},y)} \left[\text{HingeLoss}(f(\mathbf{X}), Y) \right]$

Logistic regression (log loss):

Probabilistic prediction $\hat{P}_f(y|\mathbf{x})$

$$\min_{f} \mathbb{E}_{\tilde{P}(\mathbf{x},y)} \left[\text{LogLoss}(\hat{P}_{f}(Y|\mathbf{X}), Y) \right]$$

Binary SVM and Binary Logistic Regression:

Fisher consistent: produce Bayes optimal decision in the limit

Binary SVM only:

Dual parameter sparsity

Surrogate loss for multiclass cases:

Extend binary surrogate loss like hinge-loss and log-loss to multiclass

Adversarial Prediction (Asif et. al., 2015)



Adversarial Prediction \rightarrow Optimization

Adversarial Prediction

$$\min_{\hat{P}(\hat{y}|\mathbf{x})} \max_{\check{P}(\check{y}|\mathbf{x})} \mathbb{E}_{P(\mathbf{x})\hat{P}(\hat{y}|\mathbf{x})} \left[\operatorname{loss}(\hat{Y}, \check{Y}) \right]$$

s.t.
$$\mathbb{E}_{P(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} [\phi(\mathbf{X}, \check{Y})] = \mathbb{E}_{\tilde{P}(\mathbf{x}, y)} [\phi(\mathbf{X}, Y)]$$

Minimax and Lagrangian duality

Minimization over many zero-sum games

$$\min_{\theta} \mathbb{E}_{P(\mathbf{x})} \left[\max_{\substack{\mathbf{\check{p}}\mathbf{x} \\ \mathbf{\check{p}}\mathbf{x} \\ \mathbf{\check{p}}\mathbf{x} \\ \mathbf{\check{p}}\mathbf{x} }} \min_{\substack{\mathbf{\check{p}}\mathbf{X} \\ \mathbf{\check{p}}\mathbf{x} \\ \mathbf{\check{p}}\mathbf{x} }} \mathbf{L}_{\mathbf{X},\theta}^{\prime} \mathbf{\check{p}}_{\mathbf{X}} \right]$$

where:

$$(\mathbf{L}'_{\mathbf{x},\theta})_{\hat{y},\check{y}} = \operatorname{loss}(\hat{y},\check{y}) + \theta^{\mathrm{T}}(\phi(\mathbf{x},\check{y}) - \phi(\mathbf{x},\check{y}))$$
$$\hat{\mathbf{p}}_{\mathbf{x}} = [\hat{P}(\hat{Y} = 1|\mathbf{x}) \ \hat{P}(\hat{Y} = 2|\mathbf{x}) \ \dots]^{\mathrm{T}}$$

$$\check{\mathbf{p}}_{\mathbf{x}} = [\check{P}(\check{Y} = 1 | \mathbf{x}) \; \check{P}(\check{Y} = 2 | \mathbf{x}) \; \dots]^{\mathrm{T}}$$

Example of game matrix $\mathbf{L}'_{\mathbf{x},\theta}$ for zero-one loss

$$\begin{bmatrix} \psi_{1,y}(\mathbf{x}) & \psi_{2,y}(\mathbf{x}) + 1 & \psi_{3,y}(\mathbf{x}) + 1 & \psi_{4,y}(\mathbf{x}) + 1 \\ \psi_{1,y}(\mathbf{x}) + 1 & \psi_{2,y}(\mathbf{x}) & \psi_{3,y}(\mathbf{x}) + 1 & \psi_{4,y}(\mathbf{x}) + 1 \\ \psi_{1,y}(\mathbf{x}) + 1 & \psi_{2,y}(\mathbf{x}) + 1 & \psi_{3,y}(\mathbf{x}) & \psi_{4,y}(\mathbf{x}) + 1 \\ \psi_{1,y}(\mathbf{x}) + 1 & \psi_{2,y}(\mathbf{x}) + 1 & \psi_{3,y}(\mathbf{x}) + 1 & \psi_{4,y}(\mathbf{x}) \end{bmatrix}$$
$$\psi_{j,\tilde{y}}(\mathbf{x}) = f_{j}(\mathbf{x}) - f_{\tilde{y}}(\mathbf{x}) = \theta^{\mathrm{T}} \left(\phi(\mathbf{x}, j) - \phi(\mathbf{x}, \tilde{y}) \right)$$

Inner optimization

Can be solved using Linear Programming Complexity: $\mathcal{O}(|\mathcal{Y}|^{3.5})$

Adversarial Prediction \rightarrow ERM perspective

Adversarial Prediction (optimization)

$$\min_{\theta} \mathbb{E}_{P(\mathbf{x})} \left[\max_{\mathbf{\check{p}}\mathbf{x}} \min_{\mathbf{\hat{p}}\mathbf{x}} \mathbf{\hat{p}}_{\mathbf{X}}^{\mathrm{T}} \mathbf{L}_{\mathbf{X},\theta}^{\prime} \mathbf{\check{p}}_{\mathbf{X}} \right]$$

Empirical Risk Minimization with surrogate loss:

$$\begin{split} &\operatorname{AdversarialSurrogate}(f_{\theta}(\mathbf{x}), y) = \max_{\check{\mathbf{p}}_{\mathbf{x}}} \min_{\hat{\mathbf{p}}_{\mathbf{x}}} \hat{\mathbf{p}}_{\mathbf{x}}^{\mathrm{T}} \mathbf{L}_{\mathbf{x}, \theta}' \check{\mathbf{p}}_{\mathbf{x}} \\ & \text{where: } (\mathbf{L}_{\mathbf{x}, \theta}')_{\hat{y}, \check{y}} = \operatorname{loss}(\hat{y}, \check{y}) + (f_{\theta})_{\check{y}}(\mathbf{x}) - (f_{\theta})_{(\check{y})}(\mathbf{x}) \\ & \text{Adversarial Surrogate Loss} \end{split}$$

The Nash equilibrium value of the zero-sum game characterized by matrix $L'_{x, heta}$

Outline





Adversarial Surrogate Losses for Multiclass Ordinal Classification



Ongoing and Future Works

The Adversarial Surrogate Loss for Multiclass Zero-One Classification

Based on:

Rizal Fathony, Anqi Liu, Kaiser Asif, Brian D. Ziebart. "*Adversarial Multiclass Classification: A Risk Minimization Perspective*". Advances in Neural Information Processing Systems 29 (NIPS), 2016.

Multiclass Zero-One: Related Works

Multiclass Support Vector Machine



1. The WW Model (Weston et.al., 2002)

$$\operatorname{loss}_{WW}(\mathbf{x}_i, y_i) = \sum_{j \neq y_i} [1 - (f_{y_i}(\mathbf{x}_i) - f_j(\mathbf{x}_i))]_+$$

Relative Margin Model

2. The CS Model (Crammer and Singer, 1999)

 $loss_{CS}(\mathbf{x}_i, y_i) = \max_{j \neq y_i} \left[1 - \left(f_{y_i}(\mathbf{x}_i) - f_j(\mathbf{x}_i)\right)\right]_+$ Relative Margin Model

3. The LLW Model (Lee et.al., 2004)

$$loss_{LLW}(\mathbf{x}_i, y_i) = \sum_{j \neq y_i} [1 + f_j(\mathbf{x}_i)]_+$$

with: $\sum_j f_j(\mathbf{x}_i) = 0$

Absolute Margin Model

Adversarial Prediction : Multiclass Zero-One Loss

Adversarial Game

 $\min_{\hat{P}(\hat{y}|\mathbf{x})} \max_{\check{P}(\check{y}|\mathbf{x})} \mathbb{E}_{P(\mathbf{x})\hat{P}(\hat{y}|\mathbf{x})\check{P}(\check{y}|\mathbf{x})} \left[I(\hat{Y} \neq \check{Y}) \right]$ s.t. $\mathbb{E}_{P(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} [\phi(\mathbf{X},\check{Y})] = \mathbb{E}_{\tilde{P}(\mathbf{x},y)} [\phi(\mathbf{X},Y)]$

$$\min_{\theta} \mathbb{E}_{P(\mathbf{x})} \left[\max_{\mathbf{p}_{\mathbf{x}}} \min_{\mathbf{p}_{\mathbf{x}}} \mathbf{\hat{p}}_{\mathbf{x}}^{\mathrm{T}} \mathbf{L}_{\mathbf{X},\theta}' \mathbf{\tilde{p}}_{\mathbf{x}} \right] \qquad \mathbf{L}_{\mathbf{x},\theta}' = \begin{bmatrix} \psi_{1,\tilde{y}}(\mathbf{x}) & \psi_{2,\tilde{y}}(\mathbf{x}) + 1 & \psi_{3,\tilde{y}}(\mathbf{x}) + 1 & \psi_{4,\tilde{y}}(\mathbf{x}) + 1 \\ \psi_{1,\tilde{y}}(\mathbf{x}) + 1 & \psi_{2,\tilde{y}}(\mathbf{x}) + 1 & \psi_{3,\tilde{y}}(\mathbf{x}) + 1 & \psi_{4,\tilde{y}}(\mathbf{x}) + 1 \\ \psi_{1,\tilde{y}}(\mathbf{x}) + 1 & \psi_{2,\tilde{y}}(\mathbf{x}) + 1 & \psi_{3,\tilde{y}}(\mathbf{x}) + 1 & \psi_{4,\tilde{y}}(\mathbf{x}) \end{bmatrix} \\ \psi_{j,\tilde{y}}(\mathbf{x}) = f_{j}(\mathbf{x}) - f_{\tilde{y}}(\mathbf{x}) = \theta^{\mathrm{T}} \left(\phi(\mathbf{x},j) - \phi(\mathbf{x},\tilde{y}) \right) \\ \text{shorter notation} \\ \mathbf{L}_{\mathbf{x},\theta}' = \begin{bmatrix} \psi_{1} & \psi_{2} + 1 & \psi_{3} + 1 & \psi_{4} + 1 \\ \psi_{1} + 1 & \psi_{2} & \psi_{3} + 1 & \psi_{4} + 1 \\ \psi_{1} + 1 & \psi_{2} + 1 & \psi_{3} & \psi_{4} + 1 \\ \psi_{1} + 1 & \psi_{2} + 1 & \psi_{3} + 1 & \psi_{4} \end{bmatrix} \\ \text{Nash Equilibrium} \end{cases}$$

Adversarial Zero-Sum Game (Zero-One Loss)

The augmented game for 4 classes

AL⁰⁻¹ (Adversarial Surrogate Loss) → Binary Classification





Change classification notation to $y \in \{+1, -1\}$, parameter to w and b, add L2 regularization

$$\min_{\mathbf{w},b,\boldsymbol{\xi},\boldsymbol{\delta}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_{i=1}^n (\xi_i + \delta_i)$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1 - \xi_i$
 $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge -1 - \delta_i$
 $\xi_i \ge 0, \ \delta_i \ge 0, \ i \in \{1, \dots, n\}$
Binary Al ⁰⁻¹

$$\min_{\mathbf{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1 - \xi_i$
 $\xi_i \ge 0, \ i \in \{1, \dots, n\}$

Soft Margin SVM













$AL^{0-1} \rightarrow 3$ Class Classification

AL⁰⁻¹ for 3-class zero-one classification:

Maximization over 7 hyperplanes:





$AL^{0-1} \rightarrow$ Fisher Consistency \rightarrow Property of the Minimizer



$AL^{0-1} \rightarrow Fisher Consistency$

Finding the minimizer f^*

$$\mathbf{f}^{*} = \underset{\mathbf{f}}{\operatorname{argmin}} \mathbb{E}_{P(y|\mathbf{x})} \left[\operatorname{AL}_{\mathbf{f}}^{0.1}(\mathbf{X}, Y) | \mathbf{X} = \mathbf{x} \right] = \underset{\mathbf{f}}{\operatorname{argmin}} \sum_{y=1}^{|\mathcal{Y}|} P_{y}(\mathbf{x}) \underset{S \subseteq \{1, \dots, |\mathcal{Y}|\}}{\max} \frac{\sum_{j \in S} \psi_{j,y}(\mathbf{x}) + |\mathcal{S}| - 1}{|\mathcal{S}|} \right]$$

based on the properties of the minimizer
$$\mathbf{f}^{*} = \underset{\mathbf{f}}{\operatorname{argmin}} \sum_{y=1}^{|\mathcal{Y}|} P_{y}(\mathbf{x}) \left[\frac{\sum_{j=1}^{|\mathcal{Y}|} (f_{j}(\mathbf{x}) - f_{y}(\mathbf{x})) + |\mathcal{Y}| - 1}{|\mathcal{Y}|} \right]$$

subject to $-\frac{1}{|\mathcal{Y}|} \leq f_{j}(\mathbf{x}) \leq \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|}, \quad j \in \{1, \dots, |\mathcal{Y}|\}; \quad \sum_{j=1}^{|\mathcal{Y}|} f_{j}(\mathbf{x}) = 0$
$$\mathbf{f}^{*} = \underset{\mathbf{f}}{\max} \sum_{y=1}^{|\mathcal{Y}|} P_{y}(\mathbf{x}) f_{y}(\mathbf{x})$$

subject to $-\frac{1}{|\mathcal{Y}|} \leq f_{j}(\mathbf{x}) \leq \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|}, \quad j \in \{1, \dots, |\mathcal{Y}|\}; \quad \sum_{j=1}^{|\mathcal{Y}|} f_{j}(\mathbf{x}) = 0$
$$\mathbf{Fisher Consistent}$$

 (\mathbf{x})

$AL^{0-1} \rightarrow Optimization \rightarrow Primal$

Optimization of AL⁰⁻¹ (Empirical Risk Minimization)

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \max_{\substack{\mathcal{S} \subseteq \{1, \dots, |\mathcal{Y}|\}\\\mathcal{S} \neq \emptyset}} \frac{\sum_{j \in \mathcal{S}} \theta^{\mathrm{T}} \left(\phi(\mathbf{x}_{i}, j) - \phi(\mathbf{x}_{i}, y_{i}) \right) + |\mathcal{S}| - 1}{|\mathcal{S}|}$$

Gradient for a single sample x_i

Let *R* be the set that maximize AL^{0-1} for x_i , The sub-gradient for a single sample x_i includes:

$$\frac{1}{|R|} \sum_{j \in R} \left[\phi(\mathbf{x}_i, j) - \phi(\mathbf{x}_i, y_i) \right]$$

Finding the set R:

Greedy algorithm:

- 1. Compute all $\psi_j \triangleq \theta^T (\phi(\mathbf{x}_i, j) \phi(\mathbf{x}_i, y_i))$
- 2. Sort ψ_i in non-descending order
- 3. Start with empty set $R = \emptyset$
- 4. Repeat:
- 5. Incrementally add j to the set R, update the value of AL⁰⁻¹
- 6. Until adding another one decrease the value of AL⁰⁻¹

$AL^{0-1} \rightarrow Optimization \rightarrow Dual$

Primal Quadratic Programming Formulation of AL⁰⁻¹ with L2 regularization

$$\begin{split} \min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \sup_{S \subseteq \{1,\dots,|\mathcal{Y}|\}} \frac{\sum_{j \in S} \theta^{\mathrm{T}} \left(\phi(\mathbf{x}_i, j) - \phi(\mathbf{x}_i, y_i)\right) + |S| - 1}{|S|} \\ \\ \underbrace{ \text{Constrained Primal QP} \\ \min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to:} \quad \xi_i \geq \Delta_{i,k} \quad \forall i \in \{1,\dots,n\} k \in \{1,\dots,2^{|\mathcal{Y}|} - 1\} \\ \\ \underbrace{ \text{Dual QP Formulation} \\ \max_{i=1}^n \sum_{k=1}^{n-2} \sum_{k=1}^{2^{|\mathcal{Y}|} - 1} \nu_{i,k} \alpha_{i,k} - \frac{1}{2} \sum_{i,j=1}^m \sum_{k,l=1}^{2^{|\mathcal{Y}|} - 1} \alpha_{i,k} \alpha_{j,l} \left[\Lambda_{i,k} \cdot \Lambda_{j,l} \right] \\ \\ \text{subject to} \quad \alpha_{i,k} \geq 0, \quad \sum_{k=1}^{2^{|\mathcal{Y}|} - 1} \alpha_{i,k} = C, \quad i \in \{1,\dots,n\}, \quad k \in \{1,\dots,2^{|\mathcal{Y}|} - 1\} \\ \\ \text{where:} \\ \Lambda_{i,k} = \frac{d\Delta_{i,k}}{d\theta} \text{ , and } \nu_{i,k} \text{ is the constant part of } \Delta_{i,k} \\ \\ \end{aligned}$$

$AL^{0-1} \rightarrow Optimization \rightarrow Dual \rightarrow Kernel Trick$

Kernel trick

input space rich feature space x_i $\omega(x_i)$ Compute the dot products implicitly $K(\mathbf{x}_i, \mathbf{x}_j) = \omega(\mathbf{x}_i) \cdot \omega(\mathbf{x}_j)$

Dual QP Formulation

$$\begin{split} \max_{\mathbf{\alpha}} \sum_{i=1}^{n} \sum_{k=1}^{2^{|\mathcal{Y}|}-1} \nu_{i,k} \, \alpha_{i,k} - \frac{1}{2} \sum_{i,j=1}^{m} \sum_{k,l=1}^{2^{|\mathcal{Y}|}-1} \alpha_{i,k} \alpha_{j,l} \, \underline{[\Lambda_{i,k} \cdot \Lambda_{j,l}]} \\ \text{subject to} \quad \alpha_{i,k} \ge 0, \, \sum_{k=1}^{2^{|\mathcal{Y}|}-1} \alpha_{i,k} = C, \, i \in \{1, \dots, n\}, \, k \in \{1, \dots, 2^{|\mathcal{Y}|} - 1\} \\ \Lambda_{i,k} \cdot \Lambda_{j,l} = c_{(i,k),(j,l)} \, K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{where:} \ c_{(i,k),(j,l)} = \sum_{m=1}^{|\mathcal{Y}|} \left(\frac{\mathbf{1}(m \in R_{i,k})}{|R_{i,k}|} - \mathbf{1}(m = y_i) \right) \left(\frac{\mathbf{1}(m \in R_{j,l})}{|R_{j,l}|} - \mathbf{1}(m = y_j) \right) \end{split}$$

 $R_{i,k}$ is the set of labels included in the constraint $\Delta_{i,k}$

$AL^{0-1} \rightarrow Optimization \rightarrow Dual \rightarrow Constraint Generation$

Primal and Dual Optimization

Exponential number of constraints (in primal) and dual variables

Constraint Generation Algorithm

Algorithm 1 Constraint generation method 1: Input: Training data $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n), C, \epsilon$ 2: $\theta \leftarrow \mathbf{0}$ 3: $A_i^* \leftarrow \{\Delta_{i,k} | \Delta_{i,k} = \psi_{y_i,y_i}(\mathbf{x}_i)\} \; \forall i = 1, \dots, n \; \{\text{Actual label enforces non-negativity}\}$ 4: repeat for $i \leftarrow 1, n$ do 5: $a \leftarrow \arg \max_{k \mid \Delta_i \mid k \in A_i} \Delta_{i,k}$ {Find the most violated constraint} 6: $\xi_i \leftarrow \max_{k \mid \Delta_{i,k} \in A_i^*} \Delta_{i,k}$ {Compute the example's current loss estimate} 7: if $\Delta_{i,a} > \xi_i + \epsilon$ then 8: $A_i^* \leftarrow A_i^* \cup \{\Delta_{i,a}\}$ {Add it to the enforced constraints set} 9: $\alpha \leftarrow \text{Optimize dual over } A^* = \cup_i A^*_i$ 10:Compute θ from $\boldsymbol{\alpha}$: $\theta = -\sum_{i=1}^{n} \sum_{k \mid \Delta_{i,k} \in A_{i}^{*}} \alpha_{i,k} \Lambda_{i,k}$ 11: end if 12:end for 13:14: **until** no A_i^* has changed in the iteration

Polynomial time convergence guarantee is provided

Experiment shows better convergence rate

$AL^{0-1} \rightarrow Experiments$

Dataset properties and AL⁰⁻¹ constraints

	Dataset		Prop	erties		SVM	AL^{0-1} constraints added and active			
		#class	#train	# test	#feat.	constraints	Linear	kernel	Gauss.	kernel
(1)	iris	3	105	45	4	210	213	13	223	38
(2)	$_{\mathrm{glass}}$	6	149	65	9	745	578	125	490	252
(3)	redwine	10	1119	480	11	10071	5995	1681	3811	1783
(4)	ecoli	8	235	101	7	1645	614	117	821	130
(5)	vehicle	4	592	254	18	1776	1310	311	1201	248
(6)	segment	7	1617	693	19	9702	4410	244	4312	469
(7)	sat	7	4435	2000	36	26610	11721	1524	11860	6269
(8)	opt digits	10	3823	1797	64	34407	7932	597	10072	2315
(9)	pageblocks	5	3831	1642	10	15324	9459	427	9155	551
(10)	libras	15	252	108	90	3528	1592	389	1165	353
(11)	vertebral	3	217	93	6	434	344	78	342	86
(12)	breasttissue	6	74	32	9	370	258	65	271	145

$AL^{0-1} \rightarrow Experiments \rightarrow Results$

Results for Linear Kernel and Gaussian Kernel

The **mean (standard deviation)** of the **accuracy** Bold numbers: **best** or **not significantly worse** than the best

D		Linear	Kernel		Gaussian Kernel			
	AL^{0-1}	WW	\mathbf{CS}	LLW	AL^{0-1}	WW	\mathbf{CS}	LLW
(1)	96.3 (3.1)	96.0 (2.6)	96.3 (2.4)	79.7(5.5)	96.7 (2.4)	96.4 (2.4)	96.2 (2.3)	95.4(2.1)
(2)	62.5 (6.0)	62.2 (3.6)	62.5 (3.9)	52.8(4.6)	69.5 (4.2)	66.8(4.3)	69.4 (4.8)	69.2 (4.4)
(3)	58.8 (2.0)	59.1 (1.9)	56.6(2.0)	57.7(1.7)	63.3(1.8)	64.2(2.0)	64.2(1.9)	64.7 (2.1)
(4)	86.2 (2.2)	85.7(2.5)	85.8 (2.3)	74.1(3.3)	86.0 (2.7)	84.9(2.4)	85.6 (2.4)	86.0 (2.5)
(5)	78.8 (2.2)	78.8 (1.7)	78.4 (2.3)	69.8(3.7)	84.3 (2.5)	84.4 (2.6)	83.8(2.3)	84.4 (2.6)
(6)	94.9(0.7)	94.9(0.8)	95.2 (0.8)	75.8(1.5)	96.5 (0.6)	96.6 (0.5)	96.3(0.6)	96.4(0.5)
(7)	84.9(0.7)	85.4(0.7)	84.7(0.7)	74.9(0.9)	91.9 (0.5)	92.0 (0.6)	91.9 (0.5)	91.9 (0.4)
(8)	96.6 (0.6)	96.5(0.7)	96.3(0.6)	76.2(2.2)	98.7(0.4)	98.8(0.4)	98.8(0.3)	98.9 (0.3)
(9)	96.0(0.5)	$96.1 \ (0.5)$	96.3 (0.5)	92.5(0.8)	96.8 (0.5)	96.6(0.4)	96.7(0.4)	96.6(0.4)
(10)	74.1(3.3)	72.0(3.8)	71.3(4.3)	34.0(6.4)	83.6(3.8)	83.8(3.4)	85.0 (3.9)	83.2 (4.2)
(11)	85.5 (2.9)	85.9 (2.7)	85.4 (3.3)	79.8(5.6)	86.0 (3.1)	85.3 (2.9)	85.5(3.3)	84.4(2.7)
(12)	64.4 (7.1)	59.7(7.8)	66.3 (6.9)	58.3(8.1)	68.4 (8.6)	68.1 (6.5)	66.6 (8.9)	68.0 (7.2)
avg	81.59	81.02	81.25	68.80	85.14	84.82	85.00	84.93
#b	9	6	8	0	9	6	6	7

Multiclass Zero-One Classification



Adversarial Surrogate Losses for Multiclass Ordinal Classification

Based on:

Rizal Fathony, Mohammad Bashiri, Brian D. Ziebart. "*Adversarial Surrogate Losses for Ordinal Regression*". Advances in Neural Information Processing Systems 30 (NIPS), 2017.

Ordinal Classification: Related Works

Support Vector Machine for Ordinal Classification

Extend hinge loss to ordinal classification

A. Threshold Methods (Sashua & Levin, 2003; Chu & Keerthi, 2005; Rennie & Srebro, 2005)

1. All Threshold (also called SVORIM)

$$\operatorname{AT}(\hat{f}, y) = \sum_{k=1}^{y-1} \delta(-(\eta_k - \hat{f})) + \sum_{k=y}^{|\mathcal{Y}|} \delta(\eta_k - \hat{f})$$

2. Immediate Threshold (also called SVOREX)

$$\operatorname{IT}(\hat{f}, y) = \delta(-(\eta_{y-1} - \hat{f})) + \delta(\eta_y - \hat{f})$$

 δ is a surrogate for binary classification, e.g. the hinge loss

B. Reduction Framework (Li & Lin, 2007)

Create |Y| - 1 weighted extended samples for each training sample,
Run binary classification with binary surrogate loss (e.g. hinge loss) on the extended samples

C. Cost Sensitive Classification Based Methods (Lin, 2008; Tu & Lin, 2010; Lin, 2014)

- 1. Cost Sensitive One-Versus-All (CSOVA)
- 2. Cost Sensitive One-Versus-One (CSOVO)
- 3. Cost Sensitive One-Sided-Regression (CSOSR)

Adversarial Surrogate Loss : Ordinal Classification

Adversarial Game

$$\begin{split} \min_{\hat{P}(\hat{y}|\mathbf{x})} \max_{\tilde{P}(\hat{y}|\mathbf{x})} \mathbb{E}_{P(\mathbf{x})\hat{P}(\hat{y}|\mathbf{x})} \tilde{\mathbb{E}}_{P(\hat{y}|\mathbf{x})} \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) \tilde{P}(\hat{y}|\mathbf{x}) [\phi(\mathbf{X}, \hat{Y})] &= \mathbb{E}_{\tilde{P}(\mathbf{x}, y)} [\phi(\mathbf{X}, Y)] \\ \downarrow \\ \min_{\theta} \mathbb{E}_{P(\mathbf{x})} \left[\max_{\tilde{\mathbf{p}}_{\mathbf{X}}} \min_{\tilde{\mathbf{p}}_{\mathbf{X}}} \hat{\mathbf{p}}_{\mathbf{X}}^{\mathsf{T}} \mathbf{L}_{\mathbf{X}, \theta}' \tilde{\mathbf{p}}_{\mathbf{X}} \right] \quad \mathbf{L}_{\mathbf{x}, \theta}' &= \begin{bmatrix} f_1 - f_{\tilde{y}} & f_2 - f_{\tilde{y}} + 1 & f_3 - f_{\tilde{y}} + 2 & f_4 - f_{\tilde{y}} + 3 \\ f_1 - f_{\tilde{y}} + 1 & f_2 - f_{\tilde{y}} & f_3 - f_{\tilde{y}} + 1 & f_4 - f_{\tilde{y}} + 2 \\ f_1 - f_{\tilde{y}} + 2 & f_2 - f_{\tilde{y}} + 1 & f_3 - f_{\tilde{y}} & f_4 - f_{\tilde{y}} + 1 \\ f_1 - f_{\tilde{y}} + 3 & f_2 - f_{\tilde{y}} + 2 & f_3 - f_{\tilde{y}} + 1 & f_4 - f_{\tilde{y}} \end{bmatrix} \\ \text{where } f_j(\mathbf{x}) &= \theta^{\mathsf{T}} \phi(\mathbf{x}, j) \end{split}$$

Nash Equilibrium

$$v = \operatorname{AL}_{\theta}^{\operatorname{ord}}(\mathbf{x}, y) = \max_{\substack{j,l \in \{1,\dots,|\mathcal{Y}|\}}} \frac{f_j + f_l + j - l}{2} - f_y$$
$$= \max_j \frac{f_j + j}{2} + \max_l \frac{f_l - l}{2} - f_y$$

AL^{ord} : maximization over pairs Can be independently realized

AL^{ord} → Feature Representation



$$\phi_{th}(\mathbf{x}, y) = \begin{pmatrix} y\mathbf{x} \\ I(y \le 1) \\ I(y \le 2) \\ \vdots \\ I(y \le |\mathcal{Y}| - 1) \end{pmatrix}$$

size : $m + |\mathcal{Y}| - 1$ m is the dimension of input space

- a single shared vector of feature weights - a set of threshold



An example where thresholded regression representation are useful

Multiclass Representation

$$\phi_{mc}(\mathbf{x}, y) = \begin{pmatrix} I(y=1)\mathbf{x} \\ I(y=2)\mathbf{x} \\ I(y=3)\mathbf{x} \\ \vdots \\ I(y=|\mathcal{Y}|)\mathbf{x} \end{pmatrix}$$

size : $m|\mathcal{Y}|$

m is the dimension of input spaceclass specific feature weights

An example where multiclass representation are useful



AL^{ord} → Thresholded Regression Representation

ALord for the Thresholded Regression Representation



 AL^{ord-th} : based on averaging the threshold label predictions for potentials $w \cdot x_i + 1$ and $w \cdot x_i - 1$

$AL^{ord} \rightarrow Multiclass Representation$

ALord for the Multiclass Representation

$$AL^{\text{ord-mc}}(\mathbf{x}_i, y_i) = \max_{j,l \in \{1, \dots, |\mathcal{Y}|\}} \frac{\mathbf{w}_j \cdot \mathbf{x}_i + \mathbf{w}_l \cdot \mathbf{x}_i + j - l}{2} - \mathbf{w}_{y_i} \cdot \mathbf{x}_i$$



AL^{ord-mc} : maximization over $\frac{|\mathcal{Y}|(|\mathcal{Y}|+1)}{2}$ hyperplanes

AL^{ord} → Fisher Consistency

Fisher Consistency in Ordinal Classification

$$\mathbf{f}^* = \underset{\mathbf{f}}{\operatorname{argmin}} \mathbb{E}_{P(y|\mathbf{x})} \left[\operatorname{surrogate}_{\mathbf{f}}(\mathbf{X}, Y) | \mathbf{X} = \mathbf{x} \right] \Rightarrow \underset{j}{\operatorname{argmax}} f_j^* \subseteq \underset{j}{\operatorname{argmin}} \sum_{y} P_y \left| j - y \right|$$

 $\max f_j(\mathbf{x}) = 0$ constraint is employed to remove redundant solution

Properties of the minimizer f^*

The minimizer f^* satisfies the **loss reflective** property

examples:

[-1, 0, -1, -2] [0, -1, -2, -3] [-2, -1, 0, -1, -2] [-3, -2, -1, 0, -1, -2]



Finding the minimizer f^*

$$\mathbf{f}^* = \underset{\mathbf{f}}{\operatorname{argmin}} \mathbb{E}_{P(y|\mathbf{x})} \left[\operatorname{AL}_{\mathbf{f}}^{\operatorname{ord}}(\mathbf{X}, Y) | \mathbf{X} = \mathbf{x} \right] = \underset{\mathbf{f}}{\operatorname{argmin}} \sum_{y} P_{y} \left[\underset{j,l \in \{1, \dots, |\mathcal{Y}|\}}{\max} \frac{f_{j} + f_{l} + j - l}{2} - f_{y} \right]$$

based on the **loss reflective** property

Equivalent with finding j^* (the class in a loss reflective f that has 0 value)

$$j^* = \underset{j}{\operatorname{argmin}} \sum_{y=1}^{|\mathcal{Y}|} P_y |j - y|$$
 \longrightarrow Fisher Consistent

$AL^{ord} \rightarrow Optimization \rightarrow Primal$

Optimization of AL^{ord} (Empirical Risk Minimization)

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \left[\max_{j,l \in \{1,\dots,|\mathcal{Y}|\}} \frac{\theta^{\mathrm{T}} \left[\phi(\mathbf{x}_{i},j) + \phi(\mathbf{x}_{i},l)\right] + j - l}{2} - \theta^{\mathrm{T}} \left[\phi(\mathbf{x}_{i},y_{i})\right] \right]$$

Stochastic Average Gradient (SAG) (Schimidt et.al. 2013, 2015)

Average over the gradient of each example from the last iteration it was selected Requires storing the gradient of each sample.

SAG for AL^{ord-mc}

Objective:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \left[\max_{j,l \in \{1,\dots,|\mathcal{Y}|\}} \frac{\mathbf{w}_j \cdot \mathbf{x}_i + \mathbf{w}_l \cdot \mathbf{x}_i + j - l}{2} - \mathbf{w}_{y_i} \cdot \mathbf{x}_i \right]$$

Gradient for a single sample x_i :

$$j^{*}, l^{*} = \underset{j,l \in \{1,...,|\mathcal{Y}|\}}{\operatorname{argmax}} \frac{\mathbf{w}_{j} \cdot \mathbf{x}_{i} + \mathbf{w}_{l} \cdot \mathbf{x}_{i} + j - l}{2} - \mathbf{w}_{y_{i}} \cdot \mathbf{x}_{i}$$

$$\nabla_{\mathbf{w}_{j}*} = \frac{1}{2} \mathbf{x}_{i} \qquad \nabla_{\mathbf{w}_{y_{i}}} = -\mathbf{x}_{i}$$

$$\nabla_{\mathbf{w}_{j}*} = \frac{1}{2} \mathbf{x}_{i} \qquad \nabla_{\mathbf{w}_{y_{i}}} = -\mathbf{x}_{i}$$

$$\nabla_{\mathbf{w}_{l}*} = \frac{1}{2} \mathbf{x}_{i} \qquad \nabla_{\mathbf{w}_{k}} = \mathbf{0} \quad k \in \{1, \ldots, |\mathcal{Y}|\} \setminus \{j^{*}, l^{*}, y_{i}\}$$
Store j^{*} and l^{*}
instead of full gradient

$AL^{ord} \rightarrow Optimization \rightarrow Dual$

Primal Quadratic Programming Formulation of AL^{ord} with L2 regularization

$$\min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \left[\max_{j \in 1, \dots, |\mathcal{Y}|} \frac{\theta \cdot \phi(\mathbf{x}_i, j) + j}{2} + \max_{j \in 1, \dots, |\mathcal{Y}|} \frac{\theta \cdot \phi(\mathbf{x}_i, j) - j}{2} - \theta \cdot \phi(\mathbf{x}_i, y_i) \right]$$

Constrained Primal QP

$$\min_{\theta} \frac{1}{2} \|\theta\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i + \frac{C}{2} \sum_{i=1}^n \delta_i$$

subject to: $\xi_i \ge \theta \cdot \phi(\mathbf{x}_i, j) - \theta \cdot \phi(\mathbf{x}_i, y_i) + j$ $\forall i \in \{1, \dots, n\}; j \in \{1, \dots, |\mathcal{Y}|\}$
 $\delta_i \ge \theta \cdot \phi(\mathbf{x}_i, j) - \theta \cdot \phi(\mathbf{x}_i, y_i) - j$ $\forall i \in \{1, \dots, n\}; j \in \{1, \dots, |\mathcal{Y}|\}$

Dual QP Formulation

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta}} \sum_{i,j} j(\alpha_{i,j} - \beta_{i,j}) \\ - \frac{1}{2} \sum_{i,j,k,l} (\alpha_{i,j} + \beta_{i,j}) (\alpha_{k,l} + \beta_{k,l}) (\phi(\mathbf{x}_i, j) - \phi(\mathbf{x}_i, y_i)) \cdot (\phi(\mathbf{x}_k, l) - \phi(\mathbf{x}_l, y_k)) \\ \text{subject to: } \alpha_{i,j} \ge 0; \beta_{i,j} \ge 0; \sum_j \alpha_{i,j} = \frac{C}{2}; \sum_j \beta_{i,j} = \frac{C}{2}; i, k \in \left\{1, \dots, n\}; j, l \in \{1, \dots, |\mathcal{Y}|\right\}$$

Kernel trick can also be easily applied!

$AL^{ord} \rightarrow Experiments$

Dataset properties

Dataset	# class	# train	# test	# features
diabetes	5	30	13	2
pyrimidines	5 5	51	23	27
triazines	5	130	56	60
wisconsin	5	135	59	32
machinecpu	ı 10	146	63	6
autompg	10	274	118	7
boston	5	354	152	13
stocks	5	665	285	9
abalone	10	2923	1254	10
bank	10	5734	2458	8
$\operatorname{computer}$	10	5734	2458	21
$\operatorname{calhousing}$	10	14447	6193	8

$AL^{ord} \rightarrow Experiments \rightarrow Linear Kernel$

Results for Linear Kernel

Dataset	Thr	reshold-ba	sed mode	els	Multiclass-based models				
Dataset	$\mathrm{AL}^{\mathrm{ord-th}}$	$\operatorname{RED^{th}}$	AT	IT	$\mathrm{AL}^{\mathrm{ord-mc}}$	$\operatorname{RED}^{\operatorname{mc}}$	CSOSR	CSOVO	CSOVA
diabetes	0.696	0.715	0.731	0.827	0.629	0.700	0.715	0.738	0.762
pyrimidines	0.654	0.678	0.615	0.626	0.509	0.565	0.520	0.576	0.526
triazines	0.607	0.683	0.649	0.654	0.670	0.673	0.677	0.738	0.732
wisconsin	1.077	1.067	1.097	1.175	1.136	1.141	1.208	1.275	1.338
machinecpu	0.449	0.456	0.458	0.467	0.518	0.515	0.646	0.602	0.702
autompg	0.551	0.550	0.550	0.617	0.599	0.602	0.741	0.598	0.731
boston	0.316	0.304	0.306	0.298	0.311	0.311	0.353	0.294	0.363
stocks	0.324	0.317	0.315	0.324	0.168	0.175	0.204	0.147	0.213
abalone	0.551	0.547	0.546	0.571	0.521	0.520	0.545	0.558	0.556
bank	0.461	0.460	0.461	0.461	0.445	0.446	0.732	0.448	0.989
computer	0.640	0.635	0.633	0.683	0.625	0.624	0.889	0.649	1.055
calhousing	1.190	1.183	1.182	1.225	1.164	1.144	1.237	1.202	1.601
average	0.626	0.633	0.629	0.661	0.613	0.618	0.706	0.652	0.797
# bold	5	5	4	2	5	5	2	2	1

$AL^{ord} \rightarrow Experiments \rightarrow Gaussian Kernel$

Results for Gaussian Kernel

		All Threshold	Intermediat	e Threshold
		1		
Dataset	$\mathrm{AL}^{\mathrm{ord}\text{-th}}$	SVORIM	SVOREX	
diabetes	0.696	0.665	0688	
pyrimidines	0.478	0.539	0.550	
triazines	0.609	0.612	0.604	
wisconsin	1.090	1.113	1.049	
machinecpu	0.452	0.652	0.628	
autompg	0.529	0.589	0.593	
boston	0.278	0.324	0.316	
stocks	0.103	0.099	0.100	
average	0.531	0.574	0.566	
# bold	7	3	4	

Conclusion

Conclusion



Establish connections between Adversarial Prediction and ERM



Propose Adversarial Surrogate Losses:



- Align better with the original loss
- Optimizing Adversarial Loss in the ERM Framework
- = Optimizing the original loss in the Adversarial Prediction Framework



Guarantee Fisher Consistency



Enable computational efficiency for rich feature space via kernel trick and dual parameter sparsity



Perform well in practice

Ongoing and Future Works

[1.] Taxonomy-based Classification

Adversarial Surrogate Loss for Taxonomy-based Classification



Nash equilibrium:

Analyze non-zero probability strategy of the adversary

Adversary's probability:



[1.] Taxonomy-based Classification \rightarrow Algorithm

Algorithm for finding the adversarial loss

Potentials:

 $\boldsymbol{\psi} = [0.1, 1.2, 0.5, 1.5, 0.7]$

Sorted index for each non-leaf nodes:



Complexity

 $O(md^2k^2)$

m : max # of children
d : depth of the tree
k : # of class

Ref: LP complexity: $O(k^{3.5})$



[2.] Sequence Prediction with Ordinal Classification Loss

Adversarial Surrogate Loss for Sequence Prediction with Ordinal Classification Loss

Adversarial game over joint distributions: $\min_{\theta} \mathbb{E}_{P(\mathbf{x},\mathbf{y})} \left| \max_{\substack{\check{\mathbf{p}}_{\mathbf{x}}}} \min_{\hat{\mathbf{p}}_{\mathbf{x}}} \hat{\mathbf{p}}_{\mathbf{x}}^{\mathrm{T}} \mathbf{L}_{\mathbf{X},\theta}' \check{\mathbf{p}}_{\mathbf{X}} \right|$ $\min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})\in\Xi} \mathbb{E}_{P(\mathbf{x})\hat{P}(\hat{\mathbf{y}}|\mathbf{x})\check{P}(\check{\mathbf{y}}|\mathbf{x})} \left[loss(\hat{\mathbf{Y}},\check{\mathbf{Y}}) \right]$ where: $\Xi : \mathbb{E}_{P(\mathbf{x})\check{P}(\check{\mathbf{y}}|\mathbf{x})}[\Phi(\mathbf{X},\check{\mathbf{Y}})] = \tilde{\Phi},$ $\log(\hat{\mathbf{y}},\check{\mathbf{y}}) = \sum_{t=1}^{T} |\hat{y}_t - \check{y}_t| \quad \Phi(\mathbf{x},\check{\mathbf{y}}) = \sum_{t=1}^{T} \phi(\underbrace{\mathbf{x}}_{t} y_t, y_{t+1}) \\ -\operatorname{Game analysis for ordinal classification loss guarantee}$ - Algorithm - Complexity analysis Nash equilibrium: - Algorithm implementation Exists an equilibrium where only **two** strategies $v = \max_{\mathbf{y},\mathbf{y}'} \frac{1}{2} \left[\sum_{t=1}^{T} |y_t - y'_t| + \theta^T \sum_{t=1}^{T-1} \left(\phi(\mathbf{x}, \begin{vmatrix} \mathbf{Future:} \\ y_t, y_{t+1} \\ \mathbf{Formal proof}'_{t+1} \end{pmatrix} \right] - \theta^T \sum_{t=1}^{T-1} \phi(\mathbf{x}, \tilde{y}_t, \tilde{y}_{t+1}) \\ - \text{Real data experiments} \right]$ Can be solved using dynamic programming, compleExtend analysis to other additive multiclass losses (e.g. zero-one loss)

[3.] Adversarial Graphical Model

Focus on tree structures with additive loss



Adversarial Game Over Marginal Distributions

$$\begin{split} \max_{\substack{\tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j})\\\tilde{P}_{i}(\tilde{y}_{i})}} & \min_{\hat{P}_{i}(\hat{y}_{i})} \sum_{i} \mathbb{E}_{\tilde{P}_{i}(\tilde{y}_{i})} \hat{P}_{i}(\tilde{y}_{i}) | \operatorname{bss}(\hat{y}_{i},\tilde{y}_{i}) + \sum_{i,j} \mathbb{E}_{\tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j})} \theta^{T} \phi(\mathbf{x},\tilde{y}_{i},\tilde{y}_{j}) - \sum_{i,j} \theta^{T} \phi(\mathbf{x},\tilde{y}_{i},\tilde{y}_{j}) \\ & s.t. \sum_{j} \check{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j}) = \check{P}_{i}(\tilde{y}_{i}), \forall i; \sum_{i} \check{P}_{i,j} \\ & \prod_{i,j} (\tilde{y}_{i},\tilde{y}_{j}) \sum_{i} v_{i} + \sum_{i,j} \mathbb{E}_{\tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j})} \theta^{T} \phi(\mathbf{x},\tilde{y}_{i},\tilde{y}_{j}) - \sum_{i,j} \\ & \prod_{i,j} (\tilde{y}_{i},\tilde{y}_{j}) \sum_{i} v_{i} + \sum_{i,j} \mathbb{E}_{\tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j})} \theta^{T} \phi(\mathbf{x},\tilde{y}_{i},\tilde{y}_{j}) - \sum_{i,j} \\ & \sum_{i} \tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j}) \sum_{i} v_{i} + \sum_{i,j} \mathbb{E}_{\tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j})} \theta^{T} \phi(\mathbf{x},\tilde{y}_{i},\tilde{y}_{j}) - \sum_{i,j} \\ & \sum_{i} \tilde{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j}) = \check{P}_{i}(\tilde{y}_{i}), \forall i; \sum_{i} \check{P}_{i,j}(\tilde{y}_{i},\tilde{y}_{j}) - \sum_{i,j} \\ & v_{i} \leq \sum_{l} \check{P}_{i}(l) \operatorname{loss}(k,l), \forall i, k. \end{aligned}$$

[*.] Timeline

[1.] Taxonomy Classification

- Formal proof
- More efficient implementation
- Real data experiments

[2.] Adversarial Loss for Sequence Prediction

- Formal proof for ordinal classification case
- Real data experiments
- Extend analysis to other additive multiclass losses

[3.] Adversarial Graphical Models

(with focus on tree structures)

- Better ways to solve the LP
- Adversarial surrogate loss for graphical models
- Real data experiments





Thank You