

# Adversarial Multiclass Classification: A Risk Minimization Perspective

Rizal Fathony, Anqi Liu, Kaiser Asif, and Brian D. Ziebart  
 {rfatho2, aliu33, kasif2, bziebart}@uic.edu

Department of Computer Science, University of Illinois at Chicago



## Overview

### Empirical Risk Minimization (ERM)

- ▶ Goal of classification: minimize classification error - zero-one loss
- ▶ Minimize zero-one loss over training data: NP-hard in general

### Support Vector Machine (SVM)

- ▶ Optimize hinge loss (a convex surrogate loss) over training data
- ▶ Binary SVM: Fisher consistent and universally consistent
- ▶ Generalizing SVM to multiclass case is challenging: loses consistency guarantees or does not perform well in practice

### Adversarial Classification

- ▶ Optimizes **exact loss** (zero-one) and **approximates training data**
- ▶ Promising empirical results for cost-sensitive and multivariate losses (Asif et al. 2015, Wang et al. 2015)

### Our Approach:

1. Recast zero-one adversarial classification from ERM perspective by analyzing the Nash equilibrium and define a new multiclass loss
2. Fill the long-standing gap in ERM methods by simultaneously:
  - (1) Guaranteeing Fisher consistency and universal consistency
  - (2) Enabling computational efficiency via kernel trick and dual sparsity
  - (3) Providing competitive performance in practice
3. Significantly improve computational efficiency  
 → no game solving using linear programming is required

## Related Work

### Multiclass Support Vector Machine

- ▶ Three main formulations:
  - 1) **WW** by Weston and Watkins (1999):  
 $\text{loss}_{\text{WW}}(\mathbf{x}_i, y_i) = \sum_{j \neq y_i} [1 - (f_{y_i}(\mathbf{x}_i) - f_j(\mathbf{x}_i))]_+$
  - 2) **CS** by Crammer and Singer (2002):  
 $\text{loss}_{\text{CS}}(\mathbf{x}_i, y_i) = \max_{j \neq y_i} [1 - (f_{y_i}(\mathbf{x}_i) - f_j(\mathbf{x}_i))]_+$
  - 3) **LLW** by Lee, Lin and Wahba (2004):  
 $\text{loss}_{\text{LLW}}(\mathbf{x}_i, y_i) = \sum_{j \neq y_i} [1 + f_j(\mathbf{x}_i)]_+$ , subject to  $\sum_j f_j(\mathbf{x}_i) = 0$
- ▶ WW and CS: relative margin, LLW: absolute margin
- ▶ Only LLW is Fisher consistent and universally consistent. (Tewari and Bartlett 2004, Liu 2007)
- ▶ LLW's use of absolute margin → often performs poorly in datasets with low dimensional features. (Doğan et al. 2016)

### Adversarial Prediction Games (Asif et al. 2015)

- ▶ Two player zero-sum games:
  - 1) Adversarial player: controls conditional label distribution  $\tilde{P}(y|\mathbf{x})$   
 → must approximate training data, but otherwise maximize expected loss
  - 2) Estimator player: controls  $\hat{P}(y|\mathbf{x})$  and seeks to minimize expected loss

▶ Formulation:

$$\min_P \max_{\tilde{P}: \mathbb{E}_{\tilde{P}(\mathbf{x}, y)}[\phi(\mathbf{X}, Y)] = \tilde{\phi}} \mathbb{E}_{\tilde{P}(\mathbf{x}, y)}[\text{loss}(\hat{Y}, Y)]$$

- ▶ Feature moments  $\tilde{\phi} = \mathbb{E}_{\tilde{P}(\mathbf{x}, y)}[\phi(\mathbf{X}, Y)]$ , are measured from training data

▶ For zero-one loss, it reduces to an optimization convex in  $\theta$ :

$$\min_{\theta} \sum_i \max_{\mathbf{p}} \min_{\hat{\mathbf{p}}} \hat{\mathbf{p}}_i^T \mathbf{L}'_{\mathbf{x}_i, \theta} \hat{\mathbf{p}}_{\mathbf{x}_i}; \quad \mathbf{L}'_{\mathbf{x}_i, \theta} = \begin{bmatrix} \psi_{1, y_i}(\mathbf{x}_i) & \cdots & \psi_{|\mathcal{Y}|, y_i}(\mathbf{x}_i) + 1 \\ \vdots & \ddots & \vdots \\ \psi_{1, y_i}(\mathbf{x}_i) + 1 & \cdots & \psi_{|\mathcal{Y}|, y_i}(\mathbf{x}_i) \end{bmatrix}$$

- ▶ Inner zero-sum game can be solved using a linear program

## Risk Minimization Perspective

**Theorem** The model parameters  $\theta$  for multiclass zero-one adversarial classification are equivalently obtained from empirical risk minimization under the adversarial zero-one loss function:

$$AL_{\mathbf{r}}^{0-1}(\mathbf{x}_i, y_i) = \max_{S \subseteq \{1, \dots, |\mathcal{Y}|\}, S \neq \emptyset} \frac{\sum_{j \in S} \psi_{j, y_i}(\mathbf{x}_i) + |S| - 1}{|S|}, \text{ where}$$

$S$  is any non-empty member of the powerset of classes  $\{1, 2, \dots, |\mathcal{Y}|\}$

### Plots of $AL^{0-1}$ in binary and 3-class classification

- ▶  $AL^{0-1}$  is the maximum value over  $2^{|\mathcal{Y}|} - 1$  linear hyperplanes
- ▶ Binary classification: similar with hinge loss, but with two hinges at -1 and 1 (as shown in figure on the right)
- ▶ Three class classification: the loss function has seven facets.
- ▶ Comparison with WW and CS surrogates (as shown below)

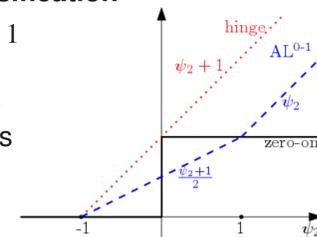


Figure: Binary  $AL^{0-1}$  over the space of potential differences  $\psi_{j,y}(\mathbf{x})$  when  $y = 1$

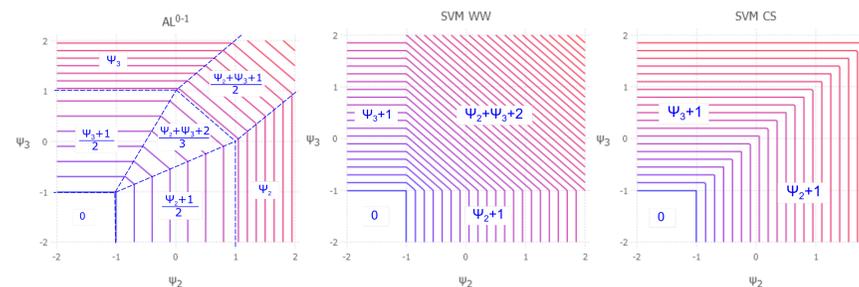


Figure: Loss function contour plots over the space of potential differences for the prediction task with three classes when  $y = 1$  under  $AL^{0-1}$ , WW and CS.

## Statistical Consistency

### Fisher Consistency

- ▶ Minimizing a Fisher consistent loss will yield the Bayes optimal decision boundary given the true distribution,  $P(x, y)$
- ▶ Multiclass: it requires  $\text{argmax}_j f_j^*(\mathbf{x}) \subseteq \text{argmax}_j P_j(\mathbf{x})$ , where  $P_j(\mathbf{x}) \triangleq P(Y = j|\mathbf{x})$  and  $\mathbf{f}^*(\mathbf{x})$  is the minimizer of  $\mathbb{E}[\text{loss}_{\mathbf{r}}(\mathbf{X}, Y)|\mathbf{X} = \mathbf{x}]$
- ▶ The minimizer of  $\mathbb{E}[AL_{\mathbf{r}}^{0-1}(\mathbf{X}, Y)|\mathbf{X} = \mathbf{x}]$  resides on the hyperplane defined by the complete set of labels,  $\mathcal{S} = \{1, \dots, |\mathcal{Y}|\}$
- ▶ It is equivalent with the area where  $-\frac{1}{|\mathcal{Y}|} \leq f_j(\mathbf{x}) \leq \frac{|\mathcal{Y}|-1}{|\mathcal{Y}|}, \forall j \in \{1, \dots, |\mathcal{Y}|\}$ , s.t.  $\sum_j f_j(\mathbf{x}) = 0$
- ▶ The minimization reduces to:

$$\max_{\mathbf{f}} \sum_{y=1}^{|\mathcal{Y}|} P_y(\mathbf{x}) f_y(\mathbf{x}) \text{ s.t. } -\frac{1}{|\mathcal{Y}|} \leq f_j(\mathbf{x}) \leq \frac{|\mathcal{Y}|-1}{|\mathcal{Y}|} \quad j \in \{1, \dots, |\mathcal{Y}|\}; \quad \sum_{j=1}^{|\mathcal{Y}|} f_j(\mathbf{x}) = 0$$

- ▶ The solution satisfies the requirement of Fisher consistency

### Universal Consistency

- ▶  $AL^{0-1}$  is a Lipschitz loss with constant 1, optimizing it with universal kernel and reasonably small regularization on any distribution yields Bayes optimal classifier. (Steinwart et al. 2008)

## Optimization

### Primal Optimization using Stochastic Sub-gradient Descent

- ▶ The sub-gradient of  $AL^{0-1}$  includes the mean of feature differences:

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

- ▶ The set  $R$  can be computed optimally using a greedy algorithm

### Dual Optimization using Quadratic Programming (QP)

- ▶ Constrained QP of  $AL^{0-1}$  plus L2 regularization

$$\min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad \xi_i \geq \Delta_{i,k} \quad \forall i \in \{1, \dots, n\}, k \in \{1, \dots, 2^{|\mathcal{Y}|-1}\}$$

where  $2^{|\mathcal{Y}|-1}$  possible constraints for example  $i$  are denoted as  $\Delta_{i,k}$

- ▶ Dual QP formulation

$$\max_{\alpha} \sum_{i=1}^n \sum_{k=1}^{2^{|\mathcal{Y}|-1}} \nu_{i,k} \alpha_{i,k} - \frac{1}{2} \sum_{i,j=1}^m \sum_{k,l=1}^{2^{|\mathcal{Y}|-1}} \alpha_{i,k} \alpha_{j,l} [\Lambda_{i,k} \cdot \Lambda_{j,l}]$$

$$\text{s.t.} \quad \alpha_{i,k} \geq 0, \quad \sum_{k=1}^{2^{|\mathcal{Y}|-1}} \alpha_{i,k} = C, \quad i \in \{1, \dots, n\}, k \in \{1, \dots, 2^{|\mathcal{Y}|-1}\}$$

where  $\Lambda_{i,k} \triangleq \frac{d\Delta_{i,k}}{d\theta}$ , and  $\nu_{i,k}$  is the constant part of  $\Delta_{i,k}$

### Kernel Trick

- ▶ Dual QP only depends on dot products
- ▶ Enables efficient rich feature expansion

### Constraint Generation

- ▶ The number of constraints in QP is exponential
- ▶ Constraint generation method → efficiently solve the problem
- ▶ Polynomial time convergence guarantee is provided

## Experiments and Results

Table: The mean and standard deviation of the accuracy for each model with linear and Gaussian kernel. Bold numbers indicate that the result is the best or not significantly worse than the best.

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

- ▶ Experiment using Linear Kernel
  - ▶ LLW performs poorly in all datasets
  - ▶  $AL^{0-1}$  has a slight advantage on average accuracy and number of "best"
- ▶ Experiment using Gaussian Kernel
  - ▶ Provides access to much richer feature spaces
  - ▶ Increases performance of all models, especially the LLW
  - ▶  $AL^{0-1}$  maintains a slight advantage

**Acknowledgments:** This research was supported as part of the Future of Life Institute (futureoflife.org) FLI-RFP-AI1 program, grant#2016-158710 and by NSF grant RI-#1526379.