

## Performance-Aligned Learning Algorithms using Distributionally Robustness Principle

## **Rizal Fathony**

Post-Doctoral Fellow @ Carnegie Melon University

Joint work with: Anqi Liu, Kaiser Asif, Mohammad Bashiri, Wei Xing, Sima Behpour, Xinhua Zhang, Brian Ziebart, Zico Kolter.

## Supervised Learning | Classification



## **Binary/Multiclass Classification**

#### Example: Digit Recognition



#### Performance Metric: Accuracy

accuracy
$$(\hat{y}, y) = \frac{1}{n} \sum_{i} I(\hat{y}_i = y_i)$$

#### Loss Metric: Zero-One Loss

$$loss(\hat{y}, y) = \frac{1}{n} \sum_{i} I(\hat{y}_i \neq y_i)$$

## **Ordinal Regression/Classification**

**Example: Movie Rating Prediction** 



Predicted vs Actual Label:

Distance 
$$\uparrow \rightarrow$$
 Loss  $\uparrow$ 

## **Classification with Imbalance Datasets**



#### **Confusion Matrix**

		Actual		
		Positive	Negative	
Pred.	Positive	True	False	Predicted
		Pos. $(TP)$	Pos. $(FP)$	Pos. $(PP)$
	Negative	False	True	Predicted
		Neg. $(FN)$	Neg. $(TN)$	Neg. (PN)
		Actual	Actual	All Data
		Pos. $(AP)$	Neg. $(AN)$	(ALL)

### Performance Metric: F1 Score

$$F1score(\widehat{\boldsymbol{y}}, \boldsymbol{y}) = \frac{2 \text{ TP}}{\text{AP} + \text{PP}}$$

## Learning Tasks & Loss/Performance Metric

Machine Learning Tasks	Popular Loss/Performance Metrics
Imbalance Datasets	<ul> <li>F1-Score</li> <li>Area under ROC Curve (AUC)</li> <li>Precision vs Recall</li> </ul>
Medical classification tasks	<ul> <li>Specificity</li> <li>Sensitivity</li> <li>Bookmaker Informedness</li> </ul>
Information retrieval tasks	<ul> <li>Precision@k</li> <li>Mean Average Precision (MAP)</li> <li>Discounted cumulative gain (DCG)</li> </ul>
Weighted classification tasks	- Cost-sensitive loss metric
Rating tasks	<ul><li>Cohen's kappa score</li><li>Fleiss' kappa score</li></ul>
Computational biology tasks	<ul> <li>Precision-Recall curve</li> <li>Matthews correlation coefficient (MCC)</li> </ul>



# Learning Framework

How to design a learning algorithm?

## Standard Approach for Learning Algorithms

Empirical Risk Minimization (ERM) [Vapnik, 1992]

- 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!

Since most of loss/performance metrics (e.g. Accuracy, F1-score) are discrete & non-continuous

## Surrogate Losses

ERM: prescribes the use of convex surrogate loss to avoid intractability

## Example: Binary Classification

Evaluation Metrics: Accuracy (Zero-One Loss)





## **Adversarial Prediction**

A distributionally robust learning framework

Adversarial Prediction (Asif et.al, 2015; Fathony et.al, 2018a) A Distributionally Robust Approach

Original discrete loss metric:

ERM: e.g. Logistic Regression

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

**Adversarial Prediction:** 

$$\min_{\mathcal{P}(\hat{y}|\mathbf{x})} \max_{\mathcal{Q}(\check{y}|\mathbf{x})\in\Xi} \mathbb{E}_{\tilde{P}(\mathbf{x})\mathcal{P}(\hat{y}|\mathbf{x})\mathcal{Q}(\check{y}|\mathbf{x})} \left[ loss(\hat{y},\check{y}) \right]$$
Predictor Adversary

Uncertainty set: moment matching on features

$$\Xi \triangleq \{ \mathcal{Q} \mid \mathbb{E}_{\tilde{P}(\mathbf{x})\mathcal{Q}(\check{y}|\mathbf{x})} [\phi(\mathbf{x},\check{y})] = \mathbb{E}_{\tilde{P}(\mathbf{x},y)} [\phi(\mathbf{x},y)] \}$$
Features
Features
Features

- Operates on the conditional distribution  $P(y|\mathbf{x})$  rather than  $P(\mathbf{x}, y)$ 

## **Adversarial Prediction: Dual Formulation**



### **Decomposable Metrics**

(Asif et.al, 2015; Fathony et.al, 2016, 2017, 2018a) Decomposable metrics:  $loss(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n} \sum_{i} loss(\hat{y}_i, y_i)$ 

Examples: Multiclass, Ordinal, Taxonomy-based, and Cost-sensitive Classification

$$\min_{\theta} \mathbb{E}_{\tilde{P}(\mathbf{x},y)} \left[ \max_{\mathcal{Q}(\check{y}|\mathbf{x})} \min_{\mathcal{P}(\hat{y}|\mathbf{x})} \mathbb{E}_{\mathcal{P}(\hat{y}|\mathbf{x})} \mathbb{E}_{\mathcal{P}(\hat{y}|\mathbf{x})} \left[ loss(\hat{y},\check{y}) + \theta^{\mathsf{T}} \left( \phi(\mathbf{x},\check{y}) - \phi(\mathbf{x},y) \right) \right] \right]$$

#### Simple loss metrics: Analytical solution

(e.g. Zero-One, Absolute Loss); by analyzing the equilibrium solution of zero-sum game.

#### More complex losses: Reformulation as a Linear Program

$$\begin{array}{ll} \max_{\mathbf{q},v} v + \mathbf{f}^{\mathsf{T}} \mathbf{q} - f_{y} & \text{where:} & \mathbf{L} \text{ is the loss in a matrix form, e.g:} \\ \text{s.t.: } \mathbf{L}_{(i,:)} \mathbf{q} \geq v \quad \forall i \in [k] & p_{i} = \hat{P}(\hat{Y} = i | \mathbf{x}) \\ q_{i} \geq 0 & \forall i \in [k] & q_{i} = \check{P}(\check{Y} = i | \mathbf{x}) \\ \mathbf{q}^{\mathsf{T}} \mathbf{1} = 1, & f_{i} = \theta^{\mathsf{T}} \phi(\mathbf{x}, i) \end{array} \qquad \begin{array}{l} \mathbf{L} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \\ \mathbf{L} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \\ \text{for a 4-class zero-one loss} \end{array}$$

# of variable = k + 1, where k = # of class

## Non-Decomposable Metrics



## Dual | Marginal Formulation: $\max_{\theta} \mathbb{E}_{\tilde{P}(\mathbf{x},\mathbf{y})} \left[ \min_{\mathcal{Q}(\check{\mathbf{y}}|\mathbf{x})} \max_{\mathcal{P}(\hat{\mathbf{y}}|\mathbf{x})} \sum_{k \in [0,n]} \sum_{l \in [0,n]} \frac{2}{k+l} \sum_{i} \mathcal{P}(\hat{y}_{i}=1, \sum_{i} \hat{y}_{i}=k|\mathbf{x}) \mathcal{Q}(\check{y}_{i}=1, \sum_{i} \check{y}_{i}=l|\mathbf{x}) \right]$ $-\theta^{\mathsf{T}} \sum_{i}^{n} \left[ \mathcal{Q}(\check{y}_{i} = 1 | \mathbf{x}) \phi(\mathbf{x}, \check{y}_{i}) - \phi(\mathbf{x}, y_{i}) \right]$ Size: $n^2$

### **Non-Decomposable Metric**

#### Example:

Binary Classification with F1-score metric

F1-score
$$(\hat{\mathbf{y}}, \mathbf{y}) = \frac{2 \text{ TP}}{\text{PP} + \text{AP}} = \frac{2 \sum_{i} \hat{y}_{i} y_{i}}{\sum_{i} \hat{y}_{i} + \sum_{i} y_{i}}$$

Dual: 
$$\max_{\theta} \mathbb{E}_{\tilde{P}(\mathbf{x},\mathbf{y})} \left[ \min_{\mathcal{Q}(\check{\mathbf{y}}|\mathbf{x})} \max_{\mathcal{P}(\hat{\mathbf{y}}|\mathbf{x})} \sum_{\hat{\mathbf{y}},\check{\mathbf{y}}} \mathcal{P}(\hat{\mathbf{y}}|\mathbf{x}) \mathcal{Q}(\check{\mathbf{y}}|\mathbf{x}) \left( \frac{2\sum_{i} \tilde{y}_{i} \tilde{y}_{i}}{\sum_{i} \hat{y}_{i} + \sum_{i} \check{y}_{i}} - \theta^{\mathsf{T}} \sum_{i}^{n} \left[ \phi(\mathbf{x},\check{y}_{i}) - \phi(\mathbf{x},y_{i}) \right] \right] \right]$$

Size:  $2^{n}$ 

Intractable!

## Generic Non-Decomposable Performance Metrics

Fathony & Kolter (in-submission)

More complex performance metric

$$\begin{aligned} \text{metric}(\hat{\mathbf{y}}, \mathbf{y}) &= \sum_{j} \frac{f_{j}(\text{ TP , TN , PP , AP })}{g_{j}(\text{ PP , AP })} \\ &= \sum_{j} \frac{f_{j}(\sum_{i} \hat{y}_{i}y_{i}, \sum_{i}(1-\hat{y}_{i})(1-y_{i}), \sum_{i} \hat{y}_{i}, \sum_{i} y_{i})}{g_{j}(\sum_{i} \hat{y}_{i}, \sum_{i} y_{i})} \end{aligned}$$

f is a linear function over TP and TN

### Dual | Marginal Formulation:

$\max_{\theta} \mathbb{E}_{\tilde{P}(\mathbf{x},\mathbf{y})} \left[ \min_{\mathcal{Q}(\check{\mathbf{y}} \mathbf{x})} \max_{\mathcal{P}(\hat{\mathbf{y}} \mathbf{x})} \sum_{k \in [0,n]} \sum_{l \in [0,n]} \sum_{j} \frac{1}{g_j(k,l)} f_j \left( \sum_{i} \mathcal{P}(\hat{y}_i) \right) \right]$	$\hat{y}_i = 1, \sum_i \hat{y}_i = k) \mathcal{Q}(\check{y}_i = 1, \sum_i \check{y}_i = l),$
---	--

$$\sum_{i} (\mathcal{P}(\hat{y}_{i}=0,\sum_{i}\hat{y}_{i}=k))(\mathcal{Q}(\check{y}_{i}=0,\sum_{i}\check{y}_{i}=l)),\ k,\ l) - \theta^{\intercal} \sum_{i}^{n} [\mathcal{Q}(\check{y}_{i}=1|\mathbf{x})\phi(\mathbf{x},\check{y}_{i}) - \phi(\mathbf{x},y_{i})]$$

Size:  $2n^2$ 

			Actual		
			Positive	Negative	
		Positive	True	False	Predicted
	_:		Pos. $(TP)$	Pos. $(FP)$	Pos. $(PP)$
	ted	Negative	False	True	Predicted
	$\mathbf{P}_{\mathbf{l}}$		Neg. $(FN)$	Neg. $(TN)$	Neg. $(PN)$
			Actual	Actual	All Data
			Pos. $(AP)$	Neg. $(AN)$	(ALL)

#### Cover most popular metrics:

e.g. Precision, Recall,  $F_{\beta}$ -score, Balanced Accuracy, Specificity, Sensitivity, Informednes, Kappa score, etc...

## Integration with Deep Learning Pipeline

Fathony & Kolter (in-submission)

#### **Programming Interface**

Enable programmers to easily incorporate custom performance metric into their deep learning pipeline

```
model = Chain(
   Conv((5, 5), 1=>20, relu),
   MaxPool((2,2)),
   Conv((5, 5), 20=>50, relu),
   MaxPool((2,2)),
   x -> reshape(x, :, size(x, 4)),
   Dense(4*4*50, 500),
   Dense(500, 1), vec
)
```

```
objective(x, y) = mean(
   logitbinarycrossentropy(model(x), y))
```

```
opt = ADAM(1e-3)
Flux.train!(objective, params(model),
    train_set, opt)
```

Leaning using binary cross entropy

\*) The codes are written in Julia

```
model = Chain(
   Conv((5, 5), 1=>20, relu), MaxPool((2,2)),
   Conv((5, 5), 20=>50, relu), MaxPool((2,2)),
   x -> reshape(x, :, size(x, 4)),
   Dense(4*4*50, 500), Dense(500, 1), vec
)
@metric F2Score
function define(::Type{F2Score}, C::ConfusionMatrix)
   return ((1 + 2^2) * C.tp) / (2^2 * C.ap + C.pp)
end
f2_score = F2Score()
enforce_special_case_positive!(f2_score)
objective(x, y) = ap_objective(model(x), y, f2_score)
```

Flux.train!(objective, params(model), train\_set, ADAM(1e-3))

Leaning using AP formulation for F2-metric



## Integration with Deep Learning Pipeline

Fathony & Kolter (in-submission)

#### Code examples for other performance metrics:

```
@metric GM_PrecRec  # Geometric Mean of Prec and Rec
function define(::Type{GM_PrecRec}, C::ConfusionMatrix)
    return C.tp / sqrt(C.ap * C.pp)
end
gpr = GM_PrecRec()
enforce_special_case_positive!(gpr)
```

Geometric mean of precision and recall

```
@metric Kappa
function define(::Type{Kappa}, C::ConfusionMatrix)
    pe = (C.ap * C.pp + C.an * C.pn) / C.all^2
    num = (C.tp + C.tn) / C.all - pe
    den = 1 - pe
    return num / den
end
kappa = Kappa()
enforce_special_case_positive!(kappa)
enforce_special_case_negative!(kappa)
```

Cohen's kappa score metric

NAME	FORMULA	
$F_{\beta}$ -score Geom. mean of Prec. & Recall Balanced Accuracy	$\frac{\frac{(1+\beta^2) \text{ TP}}{\beta^2 \text{ AP} + \text{ PP}}}{\frac{\text{TP}}{\sqrt{\text{ PP} \cdot \text{ AP}}}}$ $\frac{1}{2} \left(\frac{\text{TP}}{\text{ AP}} + \frac{\text{TN}}{\text{ AN}}\right)$	
Bookmaker Informedness	$\frac{\text{TP}}{\text{AP}} + \frac{\text{TN}}{\text{AN}} - 1$	
Markedness	$\frac{\text{TP}}{\text{PP}} + \frac{\text{TN}}{\text{PN}} - 1$	
Cohen's kappa score		
$(TP + TN) / ALL - (AP \cdot PP + AN \cdot PN) / ALL^{2}$		
$1-(AP \cdot PP + AN \cdot PN)/ALL^2$		

<pre>objective(x, y) = ap_objective(model(x),</pre>	y,	gpr)
<pre>Flux.train!(objective, params(model),</pre>		
<pre>train_set, ADAM(1e-3))</pre>		

#### \*) The codes are written in Julia



## Conclusion



## Conclusion

## **Adversarial Prediction Framework**



A distributionally robust learning framework with uncertainty set defined over the conditional distributions



- Align with the loss/performance metric by incorporating the metric into its learning objective
- Computationally efficient via marginalization technique



Easy to integrate with deep learning pipeline





## References

- Adversarial Cost-Sensitive Classification
   Kaiser Asif, Wei Xing, Sima Behpour, and Brian D. Ziebart.
   Conference on Uncertainty in Artificial Intelligence (UAI), 2015.
- Adversarial Multiclass Classification: A Risk Minimization Perspective Rizal Fathony, Anqi Liu, Kaiser Asif, Brian D. Ziebart. Advances in Neural Information Processing Systems 29 (NeurIPS), 2016.
- Adversarial Surrogate Losses for Ordinal Regression
   Rizal Fathony, Mohammad Bashiri, Brian D. Ziebart.
   Advances in Neural Information Processing Systems 30 (NeurIPS), 2017.
- Consistent Robust Adversarial Prediction for General Multiclass Classification
   Rizal Fathony, Kaiser Asif, Anqi Liu, Mohammad Bashiri, Wei Xing, Sima Behpour, Xinhua Zhang, Brian D. Ziebart.
   ArXiv preprint, 2018.
- **AP-Perf: Incorporating Generic Performance Metrics in Differentiable Learning** Rizal Fathony and Zico Kolter In submission, 2019.



