

# AP-Perf: Incorporating Generic Performance Metrics in Differentiable Learning

---

RIZAL FATHONY AND ZICO KOLTER <sup>(1),(2)</sup>

(1) Carnegie Mellon University

(2) Bosch Center for AI



## Highlight of our paper

A generic framework and programming tools  
that enable practitioners to easily  
align **training objective** with the **evaluation metric**

# Integration with ML Pipelines

Easily incorporates custom performance metrics into machine learning pipeline

```
class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(30, 100)
        self.fc2 = nn.Linear(100, 100)
        self.fc3 = nn.Linear(100, 1)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        return self.fc3(x).squeeze()
```

```
model = Net().to(device)
criterion = nn.BCEWithLogitsLoss().to(device)
```

```
optimizer.zero_grad()
objective = criterion(model(inputs), labels)
objective.backward()
optimizer.step()
```

Leaning using binary cross entropy

\*) Code in PyTorch

```
class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(30, 100)
        self.fc2 = nn.Linear(100, 100)
        self.fc3 = nn.Linear(100, 1)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        return self.fc3(x).squeeze()
```

```
class FBeta(PerformanceMetric):
    def __init__(self, beta):
        self.beta = beta

    def define(self, C: Confusion_Matrix):
        return ((1 + self.beta ** 2) * C.tp) \
            / ((self.beta ** 2) * C.ap + C.pp)
```

```
f2_score = FBeta(2)
f2_score.initialize()
f2_score.enforce_special_case_positive()
```

```
model = Net().to(device)
criterion = MetricLayer(f2_score).to(device)
```

```
optimizer.zero_grad()
objective = criterion(model(inputs), labels)
objective.backward()
optimizer.step()
```

Leaning using AP formulation for F2-metric

$F_\beta$  score  
definition

$$\frac{(1 + \beta^2) TP}{\beta^2 AP + PP}$$

# Motivation

---

# Evaluation Metrics

Example: Digit Recognition



Evaluation Metric:

Performance Metric: Accuracy

$$\text{Accuracy} = \frac{\# \text{ correct prediction}}{\# \text{ sample}}$$

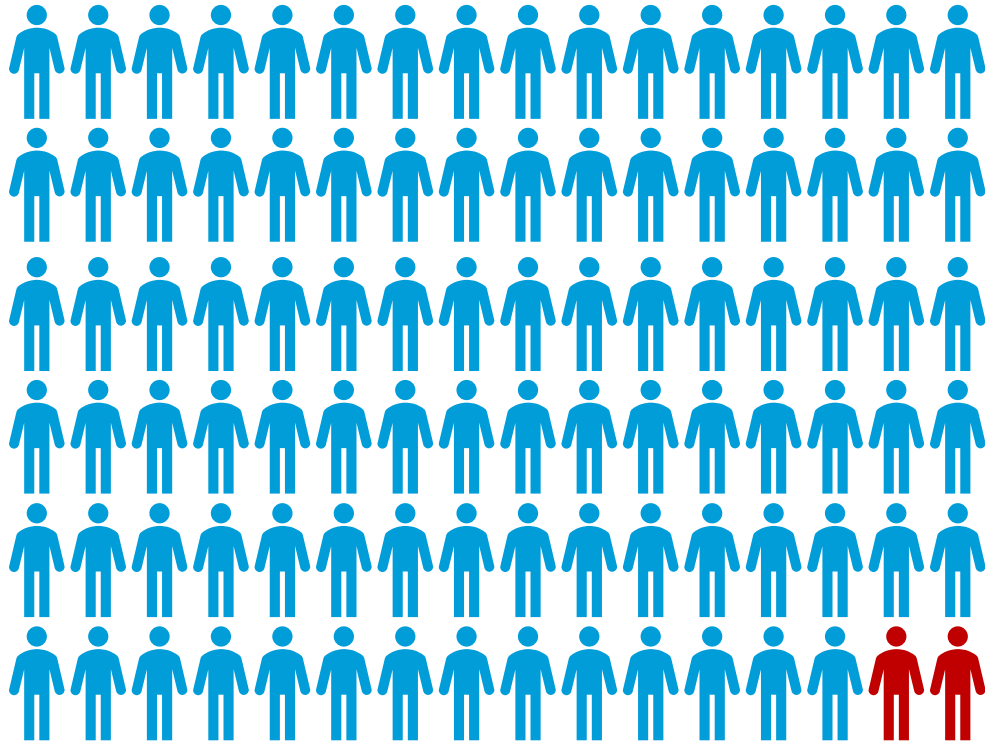
Loss Metric: Zero-One Loss

$$\text{Zero-One Loss} = \frac{\# \text{ incorrect prediction}}{\# \text{ sample}}$$

Most widely used metric!

# Accuracy metric is not always desirable

Example: Disease Prediction  
(imbalanced dataset)



Predict all samples as negative:  
Accuracy metric: 98%

## Confusion Matrix

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

$$\text{Precision} = \frac{\# \text{ true positive}}{\# \text{ predicted positive}}$$

$$\text{Recall} = \frac{\# \text{ true positive}}{\# \text{ actual positive}}$$

$$\text{Specificity} = \frac{\# \text{ true negative}}{\# \text{ actual negative}}$$

$$\text{Sensitivity} = \frac{\# \text{ true positive}}{\# \text{ actual positive}}$$

$$\text{F1-score} = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$F_{\beta}\text{-score} = \frac{(1 + \beta^2) \cdot \text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$

98% of the samples: healthy (negative samples)  
2% of the samples: have disease (positive samples)

# Learning Tasks & Evaluation Metrics

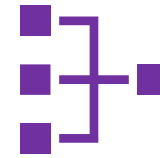
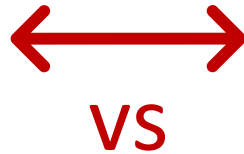
Machine Learning Tasks	Popular Evaluation Metrics
Imbalanced Datasets	<ul style="list-style-type: none"><li>- F1-Score</li><li>- Area under ROC Curve (AUC)</li><li>- Precision vs Recall</li></ul>
Medical classification tasks	<ul style="list-style-type: none"><li>- Specificity</li><li>- Sensitivity</li><li>- Bookmaker Informedness</li></ul>
Information retrieval tasks	<ul style="list-style-type: none"><li>- Precision@k</li><li>- Mean Average Precision (MAP)</li><li>- Discounted cumulative gain (DCG)</li></ul>
Weighted classification tasks	<ul style="list-style-type: none"><li>- Cost-sensitive loss metric</li></ul>
Rating tasks	<ul style="list-style-type: none"><li>- Cohen's kappa score</li><li>- Fleiss' kappa score</li></ul>
Computational biology tasks	<ul style="list-style-type: none"><li>- Precision-Recall curve</li><li>- Matthews correlation coefficient (MCC)</li></ul>

# Evaluation Metric vs Training Model Mismatch

Example: Disease prediction



Evaluation  
Metric



ML  
model



Training  
objective

Optimize specificity &  
sensitivity

Most of ML models:

- No support for specificity & sensitivity metric
- Optimize the cross-entropy objective  
(a proxy for accuracy metric)

Discrepancy:  
Evaluation metrics vs training objective



Inferior performance results  
(Cortes & Mohri, 2004; Eban et.al, 2016)





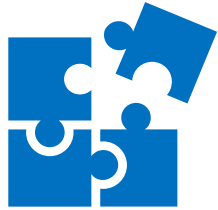
## Our paper

A generic framework and programming tools  
that enable practitioners to easily  
align **training objective** with the **evaluation metric**

# Related Works

---

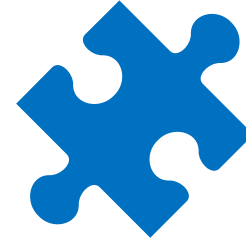
# Evaluation Metrics



## Decomposable Metrics

Can be decomposed into sample-wise sum

**Example:** accuracy, ordinal regression, and cost-sensitive metrics.



## Non-Decomposable Metrics

Cannot be decomposed into sample-wise sum

**Example:** F1-score, GPR, informedness, MCC, Kappa score.

Common in many applications

# Learning Algorithm Design

Empirical Risk Minimization Framework:

**Approximate** the **evaluation metrics** (discrete, non-continuous)  
with **convex surrogate losses**.

## Binary classification | accuracy

Evaluation Metric:

**Accuracy metric**

Convex Surrogate Losses

- ✓ **Hinge Loss** :: Support Vector Machine
- ✓ **Log Loss** :: Logistic Regression
- ✓ **Exponential Loss** :: AdaBoost

## Non-decomposable metrics

SVM-based model: **SVM-perf** (Joachims, 2005)

- ✓ Works on many complex metrics
- ✗ No statistical consistency guarantee
- ✗ Does not provide easy tool to extend the method to custom metrics

Most of other models:

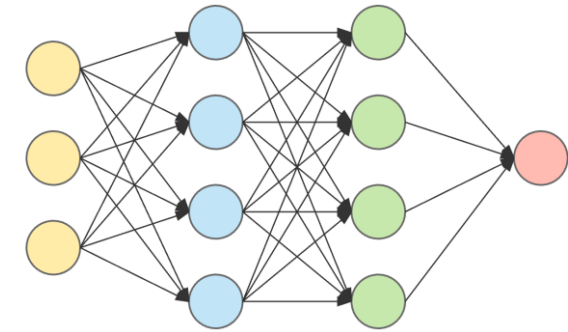
(e.g.: Koyejo et al, 2014; Narashiman et al, 2014)

- ✗ Hard to extend to custom metrics

# Neural Networks Learning

Currently the **popular** machine learning model.

Use the **classical surrogate losses** as the last layer (**objective**).



## Classification with accuracy metric

Objective: **Cross entropy objective**  
= **Logistic regression (log-loss surrogate)**

## Non-decomposable metrics

Most of **'classical' models**:

✗ Not applicable to NN learning

**NN-targeted** models:

(Eban et.al, 2016; Song et.al, 2016; Sanyal, et.al, 2018)

✗ Only support few metrics

✗ No support for custom metrics

## Practitioners' perspective



Aim to optimize an **evaluation metric** tailored **specifically** for their problem.  
(e.g. specificity, sensitivity, kappa score)



**No learning models** can easily optimize their specific evaluation metrics.

Choose the standard **cross entropy** instead

**Mismatch** between  
**Evaluation Metric vs Training Model**

# Approach

---

# Adversarial Prediction (Asif et.al, 2016; Fathony et.al, 2018)

Empirical Risk Minimization

Optimize Original  
Evaluation Metric  
Discrete, Intractable

Approximate the metric

Exact training data

Empirical Risk Minimization  
with Convex Surrogate Loss  
Convex, Tractable

More complex metric →  
Harder to construct good surrogate losses

Exact evaluation metric

Approximate training data

Adversarial Prediction  
(Asif et.al '16; Fathony et.al, '18)  
Convex, Tractable

Adversarial Prediction

No need to independently construct  
surrogate loss for every metric

# Our Method

---

OPTIMIZING GENERIC NON-DECOMPOSABLE METRICS



# Non-Decomposable Metric

Example:

Binary Classification with **F1-score** metric

$$\text{F1-score}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \text{ TP}}{\text{PP} + \text{AP}} = \frac{2 \sum_i \hat{y}_i y_i}{\sum_i \hat{y}_i + \sum_i y_i}$$

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

**AP | Decomposable metric** (Asif et.al, 2015; Fathony et.al, 2016, 2017, 2018) :

Reduces to an optimization over **sample-wise** conditional probability distributions.

$$\mathcal{P}(\hat{y}_i | \mathbf{x}_i)$$

Size: 2 (binary) x # sample

**AP | Non-decomposable metric:**

Requires optimization over **full training set** conditional probability distribution

$$\mathcal{P}(\hat{\mathbf{y}} | \mathbf{x})$$

Size:  $2^n$  (exponential)

**Marginalization technique:** optimize over marginalization distribution instead:

Original:  $\mathcal{P}(\hat{\mathbf{y}} | \mathbf{x})$   
Size:  $2^n$   
Intractable!



Marginalization:  $\mathcal{P}(\hat{y}_i = 1, \sum_i \hat{y}_i = k | \mathbf{x})$   
Size:  $n^2$  **Tractable!**

# Generic Non-Decomposable Performance Metrics

More complex performance metric

$$\text{metric}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_j \frac{a_j \text{TP} + b_j \text{TN} + f_j(\text{PP}, \text{AP})}{g_j(\text{PP}, \text{AP})}$$

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

Cover a vast range of performance metric families

Including most common use cases of non-decomposable metrics:

Precision, Recall,  $F_\beta$ -score, Balanced Accuracy, Specificity, Sensitivity, Informednes, Markedness, MCC, Kappa score, etc...

Practitioners can define their novel custom metrics

Metrics that specifically targeted to their novel problems.

Marginalization technique:

Original:  $\mathcal{P}(\hat{\mathbf{y}}|\mathbf{x})$  → Marginalization:  $\mathcal{P}(\hat{y}_i = 1, \sum_i \hat{y}_i = k|\mathbf{x})$  &  $\mathcal{P}(\hat{y}_i = 0, \sum_i \hat{y}_i = k)$

Size:  $2^n$  Size:  $2n^2$  Tractable!

Intractable!

Optimization: Gradient Descent + an ADMM-based solver (inner optimization)

# Integration with ML Pipelines

Easily incorporates custom performance metrics into machine learning pipeline

```
class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(30, 100)
        self.fc2 = nn.Linear(100, 100)
        self.fc3 = nn.Linear(100, 1)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        return self.fc3(x).squeeze()
```

```
model = Net().to(device)
criterion = nn.BCEWithLogitsLoss().to(device)
```

```
optimizer.zero_grad()
objective = criterion(model(inputs), labels)
objective.backward()
optimizer.step()
```

Leaning using binary cross entropy

\*) Code in PyTorch

```
class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(30, 100)
        self.fc2 = nn.Linear(100, 100)
        self.fc3 = nn.Linear(100, 1)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        return self.fc3(x).squeeze()
```

```
class FBeta(PerformanceMetric):
    def __init__(self, beta):
        self.beta = beta

    def define(self, C: Confusion_Matrix):
        return ((1 + self.beta ** 2) * C.tp) \
            / ((self.beta ** 2) * C.ap + C.pp)
```

```
f2_score = FBeta(2)
f2_score.initialize()
f2_score.enforce_special_case_positive()
```

```
model = Net().to(device)
criterion = MetricLayer(f2_score).to(device)
```

```
optimizer.zero_grad()
objective = criterion(model(inputs), labels)
objective.backward()
optimizer.step()
```

Leaning using AP formulation for F2-metric

$F_\beta$  score  
definition

$$\frac{(1 + \beta^2) TP}{\beta^2 AP + PP}$$

# AP-Perf: supports a wide variety of evaluation metrics

Code examples for other performance metrics:

Geometric Mean of Precision and Recall (GPR)

$$\frac{TP}{\sqrt{PP \cdot AP}}$$

```
class GM_PrecRec(PerformanceMetric):
    def define(self, C: Confusion_Matrix):
        return C.tp / sqrt(C.ap * C.pp)

gpr = GM_PrecRec()
gpr.initialize()
gpr.enforce_special_case_positive()
```

Cohen's Kappa score

$$\frac{(TP + TN) / ALL - (AP \cdot PP + AN \cdot PN) / ALL^2}{1 - (AP \cdot PP + AN \cdot PN) / ALL^2}$$

```
class Kappa(PerformanceMetric):
    def define(self, C: Confusion_Matrix):
        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

kappa = Kappa()
kappa.initialize()
kappa.enforce_special_case_positive()
kappa.enforce_special_case_negative()
```

# Novel Custom Metrics

## Write-your-own Novel Metrics

```
class NovelMetric(PerformanceMetric):  
    def define(self, C: Confusion_Matrix):  
        # write the definition of your new metric
```

### Example:

a weighted modification to the Cohen's Kappa score and the Mathews correlation coefficient (MCC)

$$0.3 \cdot \frac{(0.7 \text{ TP} + 0.3 \text{ TN}) / \text{ALL} - (0.7 \cdot \text{AP} \cdot \text{PP} + 0.3 \cdot \text{AN} \cdot \text{PN}) / \text{ALL}^2}{1 - (0.7 \cdot \text{AP} \cdot \text{PP} + 0.3 \cdot \text{AN} \cdot \text{PN}) / \text{ALL}^2} + 0.7 \cdot \frac{\text{TP} / \text{ALL} - (\text{AP} \cdot \text{PP}) / \text{ALL}^2}{\sqrt{\text{AP} \cdot \text{PP} \cdot \text{AN} \cdot \text{PN} / \text{ALL}^2}}$$

```
class NovelMetric(PerformanceMetric):  
    def define(self, C: Confusion_Matrix):  
        pe = (0.7 * C.ap * C.pp + 0.3 * C.an * C.pn) / C.all**2  
        num = (0.7 * C.tp + 0.3 * C.tn) / C.all - pe  
        den = 1 - pe  
        kappa = num / den  
  
        num2 = C.tp / C.all - (C.ap * C.pp) / C.all**2  
        den2 = ap_perf.sqrt(C.ap * C.pp * C.an * C.pn) / C.all**2  
        mcc = num2 / den2  
  
        return 0.3 * kappa + 0.7 * mcc
```

```
novel_metric = NovelMetric()  
novel_metric.initialize()  
novel_metric.enforce_special_case_positive()  
novel_metric.enforce_special_case_negative()
```

# Empirical Results

## Datasets:

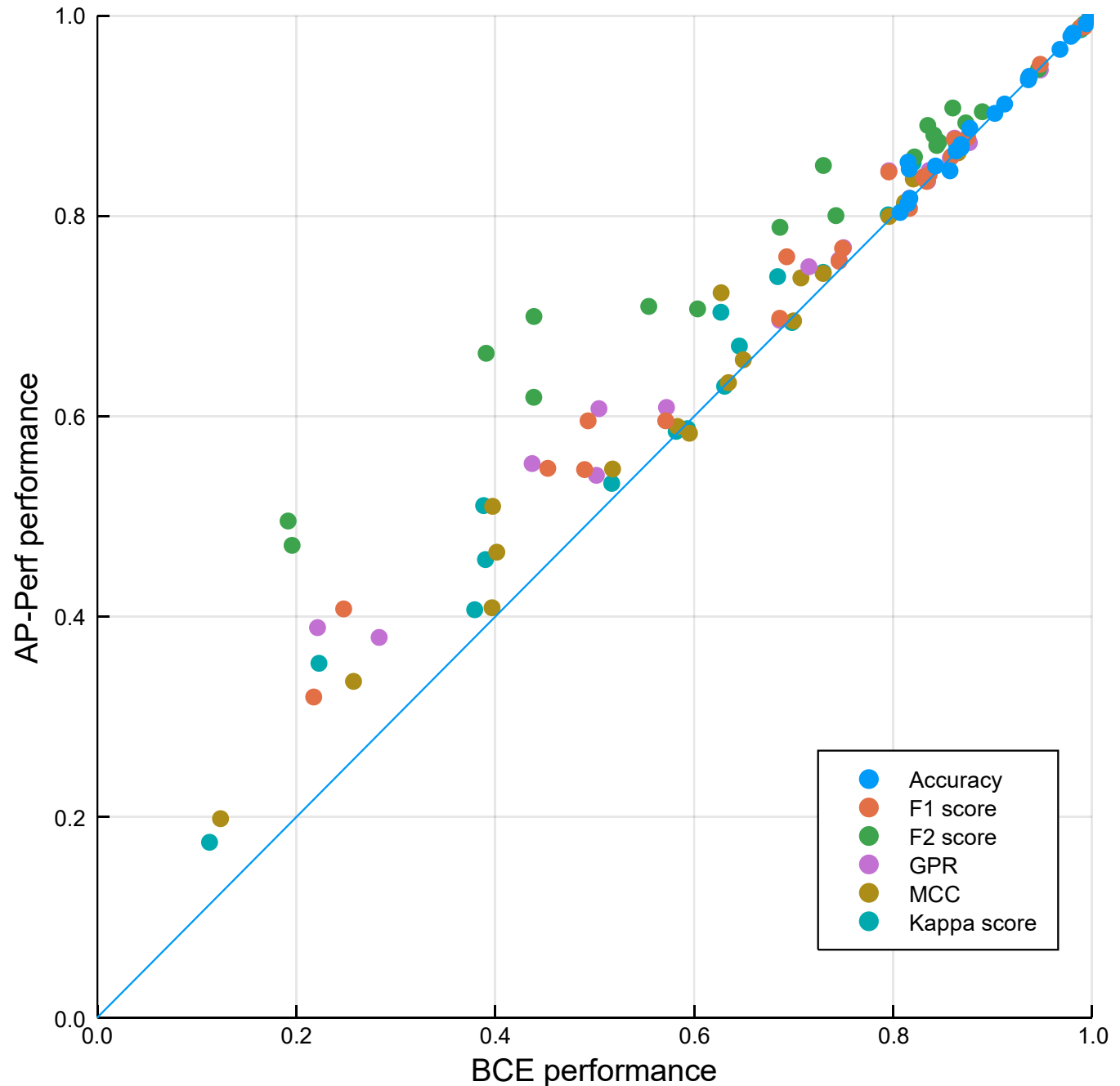
20 UCI Datasets,  
MNIST, Fashion MNIST

## Neural Networks:

Multi Layer Perceptron,  
Convolutional NN

## Performance Metrics:

- 1) Accuracy
- 2) F1 score
- 3) F2 score
- 4) Geom. Prec. Rec. (GPR)
- 5) Mathews Cor. Coef. (MCC)
- 6) Cohen's Kappa score



# Summary

	Statistical Consistency	Support Neural Network Learning	Support Custom Metrics	Easy Interface for Practitioners (to optimize custom metrics)
<b>SVM-Perf</b> (Joachim, 2005)	✗	✗	✓	✗
<b>Plug-in based classifiers</b> (Koyejo et al, 2014; Narashiman et al, 2014)	✓	✗	✗	✗
<b>Global objectives</b> (Eban et al, 2014)	✗	✓	✗	✗
<b>DAME &amp; DUPLÉ</b> (Sanyal et al, 2018)	✗	✓	✗	✗
<b>AP-Perf</b> (Fathony & Kolter, our method)	✓	✓	✓	✓

# Download



<http://proceedings.mlr.press/v108/fathony20a.html>



install:

```
pip install ap_perf
```

github:

[https://github.com/rizalzaf/ap\\_perf](https://github.com/rizalzaf/ap_perf)



install:

```
]add AdversarialPrediction
```

github:

<https://github.com/rizalzaf/AdversarialPrediction.jl>



# Thank You

---