



Being Properly Improper

Tyler Sypherd



Richard Nock

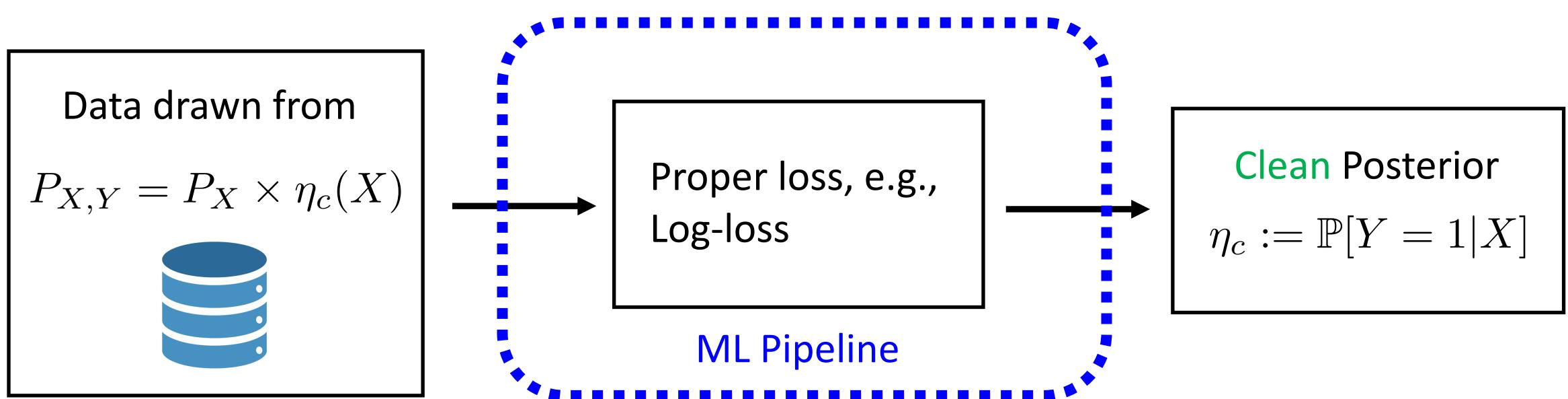


Lalitha Sankar



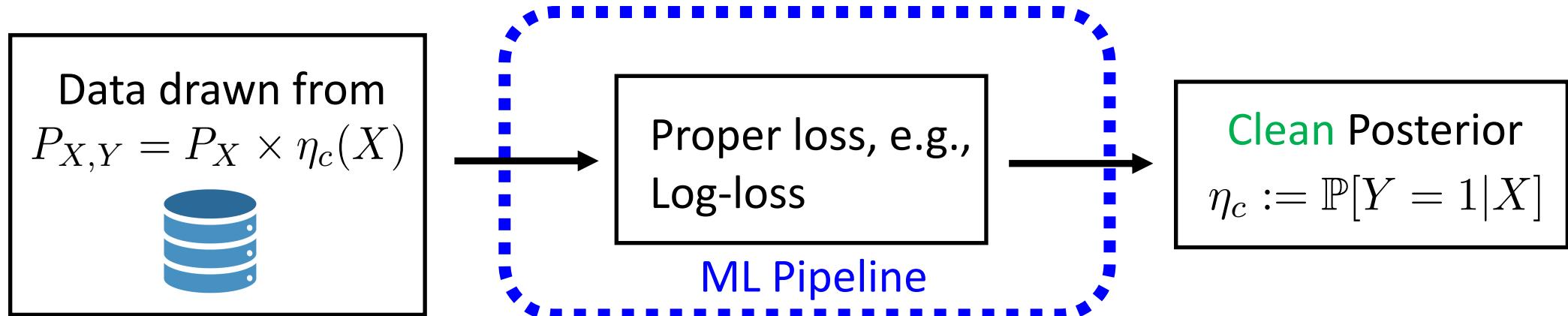
Properness

- Proper loss estimates the **true posterior**



Properness

- Proper loss estimates the **true posterior**

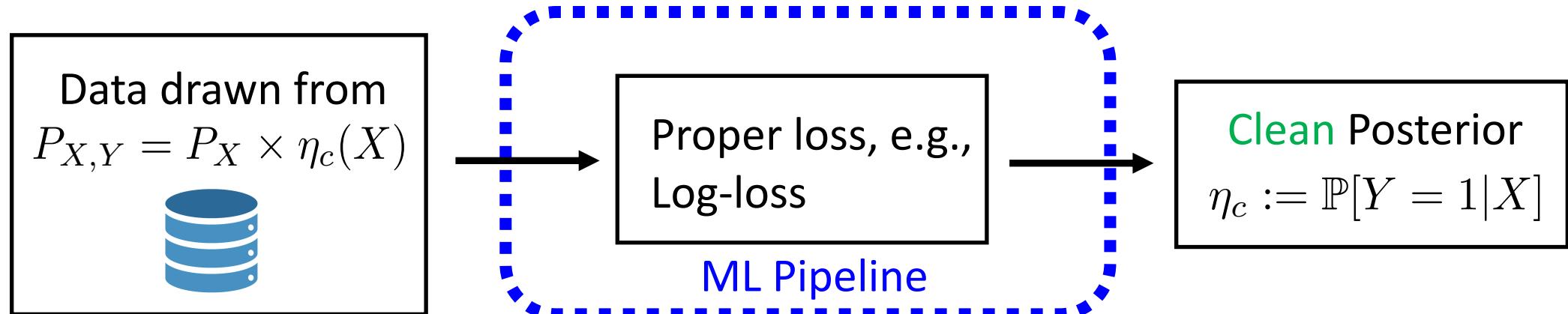


- Conditional risk of ℓ of estimated posterior $\eta \in [0, 1]$ wrt ground truth $\eta_c \in [0, 1]$

$$L(\eta, \eta_c) := \mathbb{E}_{Y \sim \text{Ber}(\eta_c)}[\ell(Y, \eta)] = \eta_c \ell_1(\eta) + (1 - \eta_c) \ell_{-1}(\eta)$$

Properness

- Proper loss estimates the **true posterior**



- Conditional risk of ℓ of estimated posterior $\eta \in [0, 1]$ wrt ground truth $\eta_c \in [0, 1]$

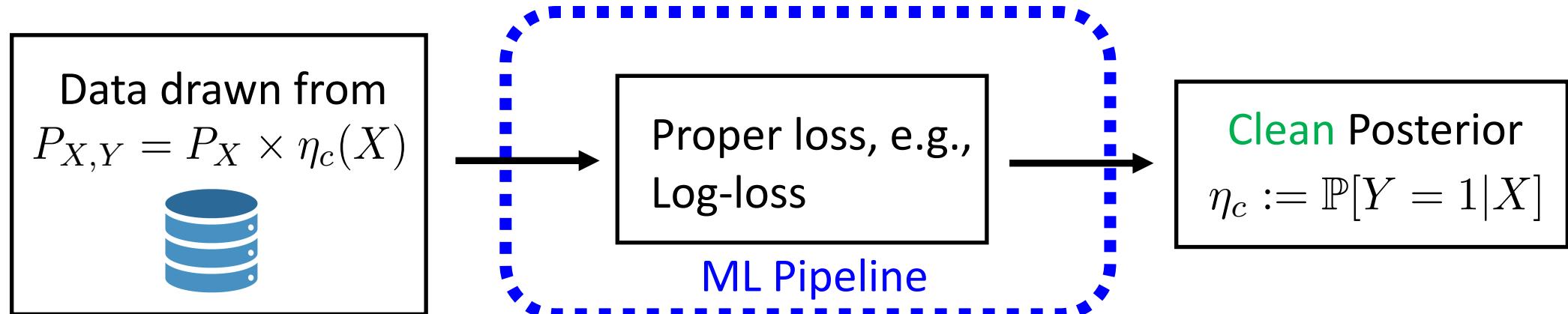
$$L(\eta, \eta_c) := \mathbb{E}_{Y \sim \text{Ber}(\eta_c)}[\ell(Y, \eta)] = \eta_c \ell_1(\eta) + (1 - \eta_c) \ell_{-1}(\eta)$$

- Proper loss:

$$\arg \min_{\eta} L(\eta, \eta_c) = \eta_c$$

Properness

- Proper loss estimates the **true posterior**



- Conditional risk of ℓ of estimated posterior $\eta \in [0, 1]$ wrt ground truth $\eta_c \in [0, 1]$

$$L(\eta, \eta_c) := \mathbb{E}_{Y \sim \text{Ber}(\eta_c)}[\ell(Y, \eta)] = \eta_c \ell_1(\eta) + (1 - \eta_c) \ell_{-1}(\eta)$$

- Example: log-loss

$$\ell_1(\eta) := -\log \eta$$

$$\ell_{-1}(\eta) := -\log(1 - \eta)$$

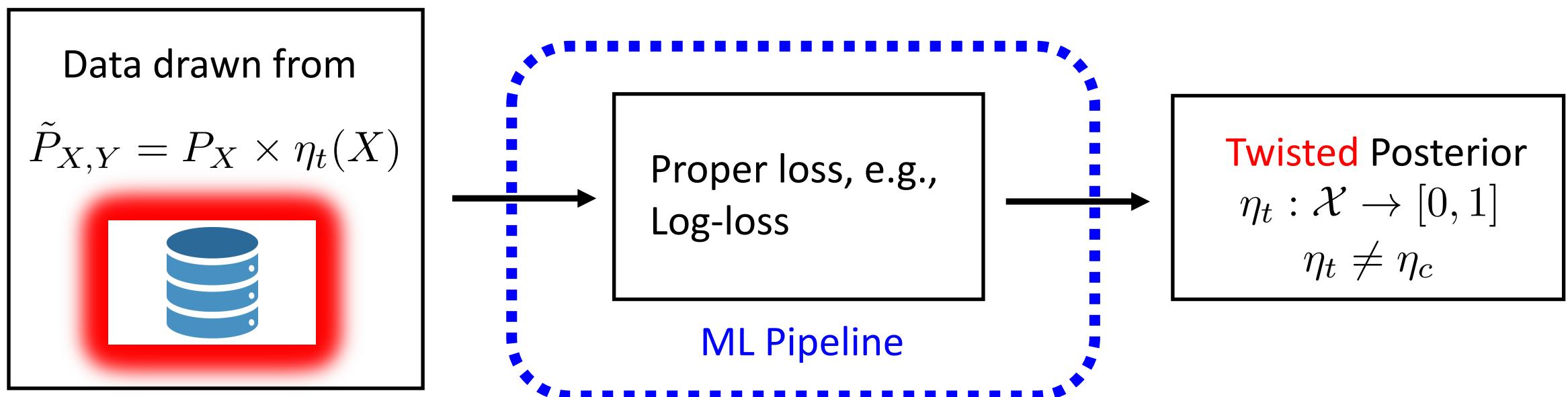
² Reid and Williamson. "Composite binary losses." (2010).

Reid and Williamson. "Information, divergence and risk for binary experiments." (2011).

Cross-entropy loss is built with log-loss

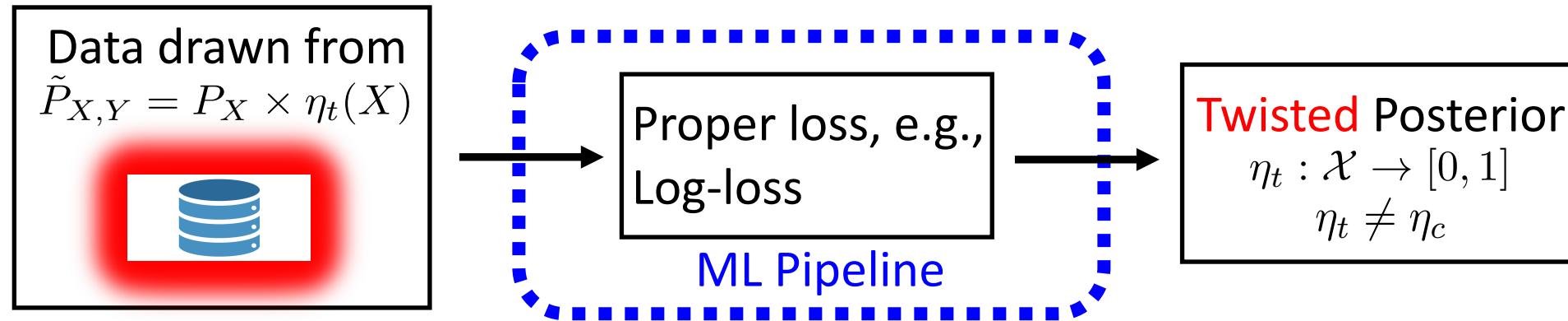
The Problem with Properness

- What happens if the data is twisted? (e.g., **label/feature/adversarial noise**)



The Problem with Properness

- What happens if the data is twisted? (e.g., **label/feature/adversarial noise**)

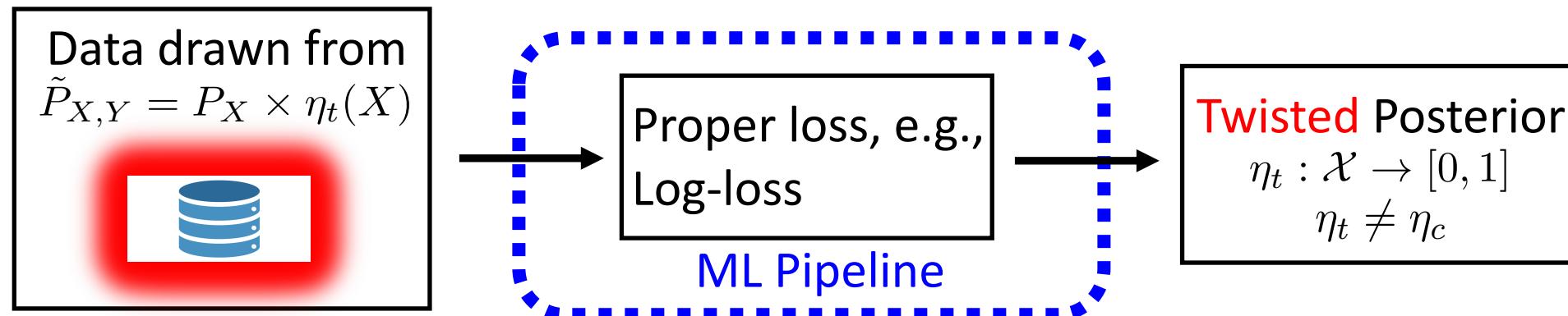


- Proper loss minimizes and returns:

$$\arg \min_{\eta} L(\eta, \eta_t) = \arg \min_{\eta} \eta_t \ell_1(\eta) + (1 - \eta_t) \ell_{-1}(\eta) = \eta_t$$

The Problem with Properness

- What happens if the data is twisted? (e.g., **label/feature/adversarial noise**)



- Proper loss minimizes and returns:

$$\arg \min_{\eta} L(\eta, \eta_t) = \arg \min_{\eta} \eta_t \ell_1(\eta) + (1 - \eta_t) \ell_{-1}(\eta) = \eta_t$$

- A **twist** $\eta_c \rightarrow \eta_t$ refers to a general mapping (label/feature/adversarial noise)

Twist-Proper Losses: Untwisting Twisted Posterior

- A loss is **twist-proper** iff for any **twist** $\eta_c \rightarrow \eta_t$, there exist hyperparameter(s) β^* , s.t.,

$$\arg \min_{\eta} L^{\beta^*}(\eta, \eta_t) = \arg \min_{\eta} \eta_t \ell_1^{\beta^*}(\eta) + (1 - \eta_t) \ell_{-1}^{\beta^*}(\eta) = \eta_c$$

- This means for every $x \in \mathcal{X}$, a twist-proper loss has hyperparameter(s) β^* that untwist the twisted posterior into the clean posterior

Twist-Proper Losses: Untwisting Twisted Posterior

- A loss is **twist-proper** iff for any **twist** $\eta_c \rightarrow \eta_t$, there exist hyperparameter(s) β^* , s.t.,

$$\arg \min_{\eta} L^{\beta^*}(\eta, \eta_t) = \arg \min_{\eta} \eta_t \ell_1^{\beta^*}(\eta) + (1 - \eta_t) \ell_{-1}^{\beta^*}(\eta) = \eta_c$$

- This means for every $x \in \mathcal{X}$, a twist-proper loss has hyperparameter(s) β^* that untwist the twisted posterior into the clean posterior
- Not vacuous: Focal loss and Super Loss are not twist-proper
- Alpha-loss (generalizes log-loss) studied by Sypherd et al. (2022) is **twist-proper**

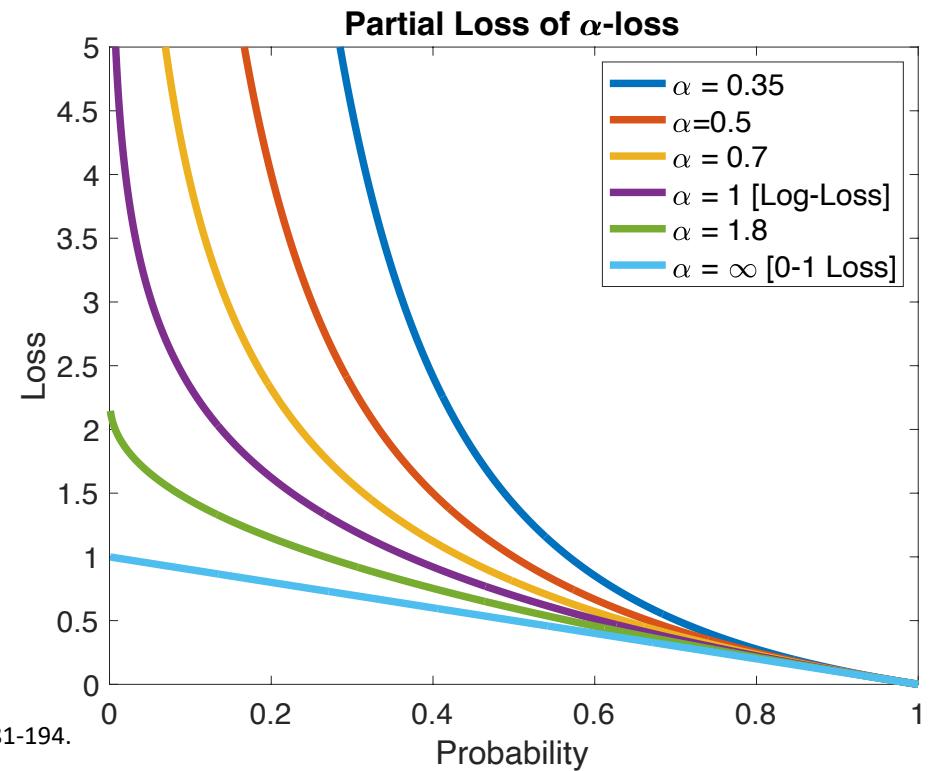
α -loss is Twist-Proper

- For $\alpha \geq 0$, partial losses $\ell_1^\alpha(\eta) := \ell_{-1}^\alpha(1 - \eta), \forall \eta \in [0, 1]$,

$$\ell_1^\alpha(\eta) := \frac{\alpha}{\alpha - 1} \left(1 - \eta^{\frac{\alpha-1}{\alpha}}\right)$$

- And, by continuity:

$$\ell_1^0(\eta) := \infty \quad \ell_1^1(\eta) := -\log \eta \quad \ell_1^\infty(\eta) := 1 - \eta$$



α -loss is Twist-Proper

- For $\alpha \geq 0$, partial losses $\ell_1^\alpha(\eta) := \ell_{-1}^\alpha(1 - \eta), \forall \eta \in [0, 1]$,

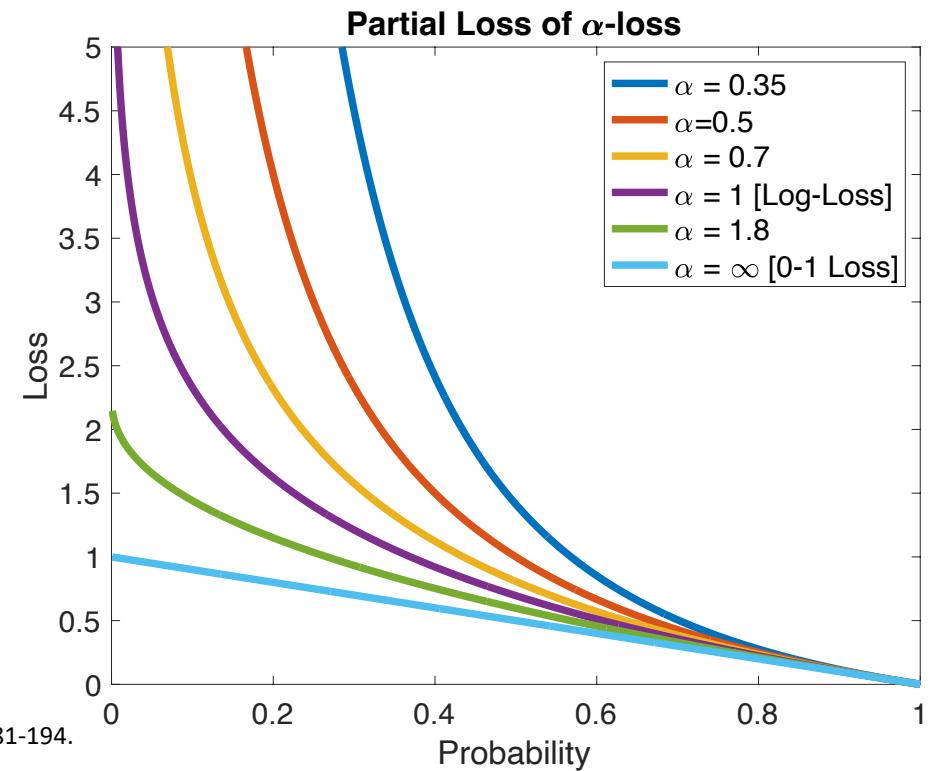
$$\ell_1^\alpha(\eta) := \frac{\alpha}{\alpha - 1} \left(1 - \eta^{\frac{\alpha-1}{\alpha}}\right)$$

- And, by continuity:

$$\ell_1^0(\eta) := \infty \quad \ell_1^1(\eta) := -\log \eta \quad \ell_1^\infty(\eta) := 1 - \eta$$

- Twist-proper (via alpha-tilted distribution):

$$\arg \min_\eta L^\alpha(\eta, \eta_t) = \frac{\eta_t^\alpha}{\eta_t^\alpha + (1 - \eta_t)^\alpha}$$



Untwisting \mathcal{X} with a fixed $\alpha_0 > 1$

Theorem: For any strictly Bayes blunting twist $\eta_c \rightarrow \eta_t$, there exists a fixed $\alpha_0 > 1$ which induces the following ordering:

$$D_{\text{KL}}(\eta_c, \eta_t; 1) > D_{\text{KL}}(\eta_c, \eta_t; \alpha_0)$$

- Strictly Bayes blunting twist $(\eta_c < \eta_t < 1/2) \vee (\eta_c > \eta_t > 1/2)$

Untwisting \mathcal{X} with a fixed $\alpha_0 > 1$

Theorem: For any strictly Bayes blunting twist $\eta_c \rightarrow \eta_t$, there exists a fixed $\alpha_0 > 1$ which induces the following ordering:

$$D_{\text{KL}}(\eta_c, \eta_t; 1) > D_{\text{KL}}(\eta_c, \eta_t; \alpha_0)$$

- Strictly Bayes blunting twist $(\eta_c < \eta_t < 1/2) \vee (\eta_c > \eta_t > 1/2)$
- This answers **yes to untwisting** the feature space with a fixed $\alpha_0 > 1$
- SLN $\eta_t = p(1 - \eta_c) + (1 - p)\eta_c$ is a strictly Bayes blunting twist for $0 < p < 1/2$

Pseudo-Inverse-Link Boost

- α -loss has been investigated in logistic regression and neural networks
- Novel boosting algorithm in order to boost hyperparameterized loss

Algorithm 1 PILBOOST

Input sample $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$, number of iterations T , $a_f > 0$, PIL \tilde{f} ;

Step 1 : let $\beta \leftarrow \mathbf{0}$; // first classifier, $H_0 = 0$

Step 2 : **for** $t = 1, 2, \dots, T$

 Step 2.1 : **for** $i = 1, 2, \dots, m$, let $w_i \leftarrow \tilde{f}(-y_i H_\beta(x_i))$ // PIL weights

 Step 2.2 : let $j \leftarrow \text{WL}(\mathcal{S}, \mathbf{w})$

 Step 2.3 : let $\eta_j \leftarrow (1/m) \cdot \sum_i w_i y_i h_j(\mathbf{x}_i)$

 Step 2.4 : let $\beta_j \leftarrow \beta_j + a_f \eta_j$

Return H_β .

Pseudo-Inverse-Link Boost

- α -loss has been investigated in logistic regression and neural networks
- Novel boosting algorithm in order to boost hyperparameterized loss

Algorithm 1 PILBOOST

Input sample $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$, number of iterations T , $a_f > 0$, PIL \tilde{f} ;
Step 1 : let $\beta \leftarrow \mathbf{0}$; // first classifier, $H_0 = 0$

Step 2 : **for** $t = 1, 2, \dots, T$

 Step 2.1 : **for** $i = 1, 2, \dots, m$, let $w_i \leftarrow \tilde{f}(-y_i H_\beta(\mathbf{x}_i))$ // PIL weights

 Step 2.2 : let $j \leftarrow \text{WL}(\mathcal{S}, \mathbf{w})$

 Step 2.3 : let $\eta_j \leftarrow (1/m) \cdot \sum_i w_i y_i h_j(\mathbf{x}_i)$

 Step 2.4 : let $\beta_j \leftarrow \beta_j + a_f \eta_j$

Return H_β .

- PILBoost yields boosting compliant rates
- Gains over AdaBoost and XGBoost on twisted data (label/feature/adversarial noise)

Pseudo-Inverse-Link Boost

- α -loss has been investigated in logistic regression and neural networks
- Novel boosting algorithm in order to boost hyperparameterized loss

Algorithm 1 PILBOOST

Input sample $\mathcal{S} = \{(\mathbf{x}_i, y_i), i = 1, 2, \dots, m\}$, number of iterations T , $a_f > 0$, PIL \tilde{f} ;

Step 1 : let $\beta \leftarrow \mathbf{0}$; // first classifier, $H_0 = 0$

Step 2 : **for** $t = 1, 2, \dots, T$

 Step 2.1 : **for** $i = 1, 2, \dots, m$, let $w_i \leftarrow \tilde{f}(-y_i H_\beta(\mathbf{x}_i))$ // PIL weights

 Step 2.2 : let $j \leftarrow \text{WL}(\mathcal{S}, \mathbf{w})$

 Step 2.3 : let $\eta_j \leftarrow (1/m) \cdot \sum_i w_i y_i h_j(\mathbf{x}_i)$

 Step 2.4 : let $\beta_j \leftarrow \beta_j + a_f \eta_j$

Return H_β .

- PILBoost yields boosting compliant rates
- Gains over AdaBoost and XGBoost on twisted data (label/feature/adversarial noise)