

# GAPTRON:

## Exploiting the Surrogate Gap in Online Multiclass Classification



Full information setting



Bandit setting

Dirk van der Hoeven



Universiteit  
Leiden  
The Netherlands

### Setting: Online Multiclass Classification

The online multiclass classification setting proceeds in rounds  $t = 1, \dots, T$ . In each round  $t$

1 the environment picks an outcome  $y_t \in \{1, \dots, K\}$  and reveals a feature vector  $\mathbf{x}_t$  to the learner

2 the learner issues a (randomized) prediction  $\hat{y}_t$

3 **Full Information Setting:** the environment reveals true outcome  $y_t$   
**Bandit setting:** the environment reveals loss  $\mathbb{1}[y_t \neq \hat{y}_t]$

4 the learner suffers  $\mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]]$

**Goal:** minimize the expected surrogate regret  $\mathcal{R}_T$

$$\mathcal{R}_T = \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}[y_t \neq \hat{y}_t] - \underbrace{\ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t)}_{\text{margin}} \right]$$

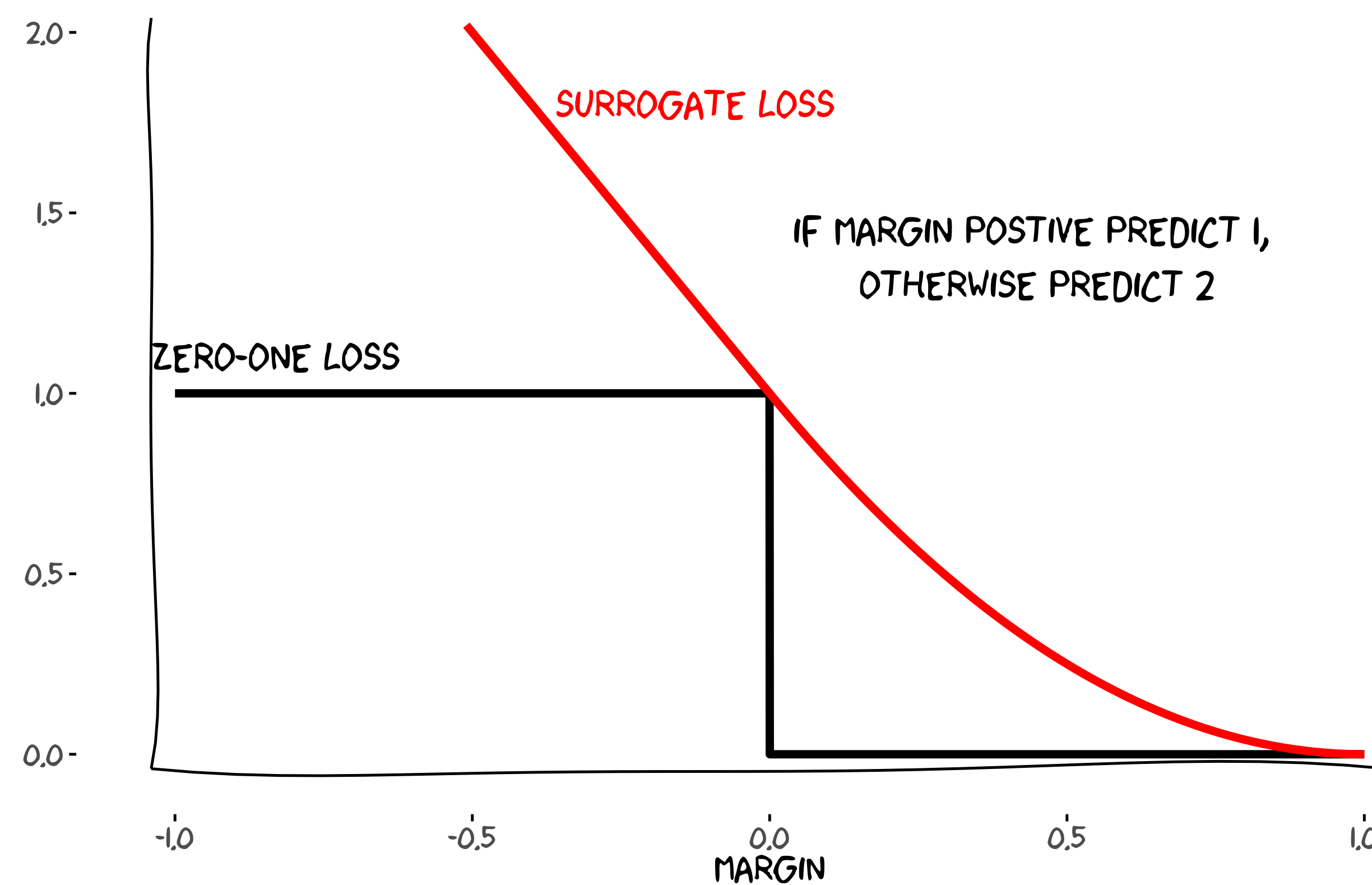
### Results

Algorithm	Full information $\mathcal{R}_T$	Bandit $\mathcal{R}_T$	Time (per round)
Standard first-order	$O(\ \mathbf{U}\ \sqrt{T})$	$O((K)^{1/3}T^{2/3})$	$O(dK)$
Standard second-order	$O(e^{\ \mathbf{U}\ } dK \ln(T))$	$O(K\sqrt{dT \ln(T)})$	$O((dK)^2)$
Gaptron	$O(K\ \mathbf{U}\ ^2)$	$O(K\sqrt{T})$	$O(dK)$

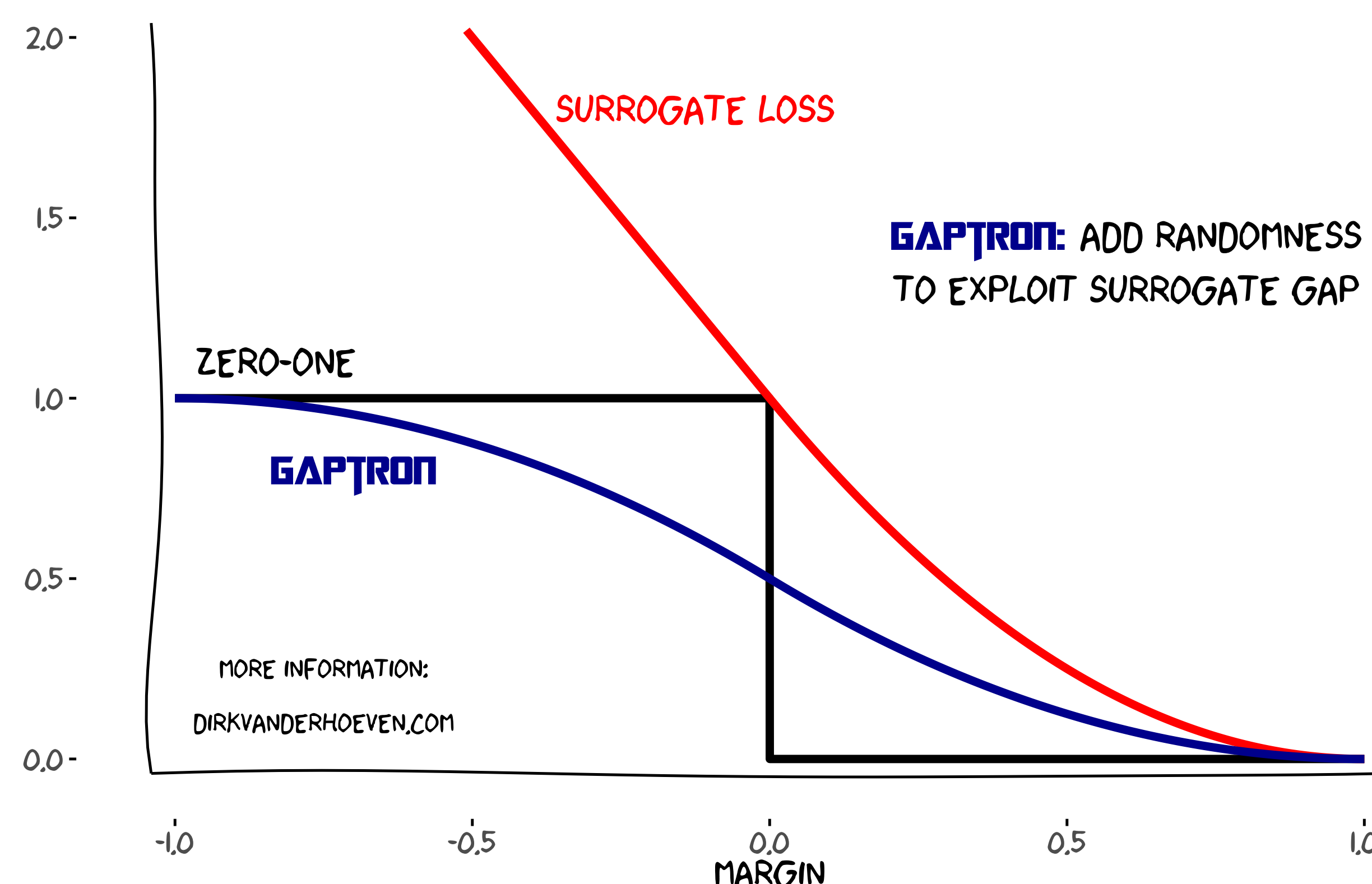
### Standard Analysis

$$\begin{aligned} & \sum_{t=1}^T \mathbb{1}[y_t \neq \hat{y}_t] - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t) \\ &= \underbrace{\left( \sum_{t=1}^T \mathbb{1}[y_t \neq \hat{y}_t] - \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) \right)}_{\text{Very wasteful: bound by 0}} + \underbrace{\left( \sum_{t=1}^T \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t) \right)}_{\text{controlled by OGD}} \\ &\leq \underbrace{\sum_{t=1}^T \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t)}_{\text{controlled by OGD}} \leq \frac{\|\mathbf{U}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla_{\mathbf{W}_t} \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t)\|^2 \\ &= O(\|\mathbf{U}\|\sqrt{T}) \end{aligned}$$

### Standard Analysis: Figure

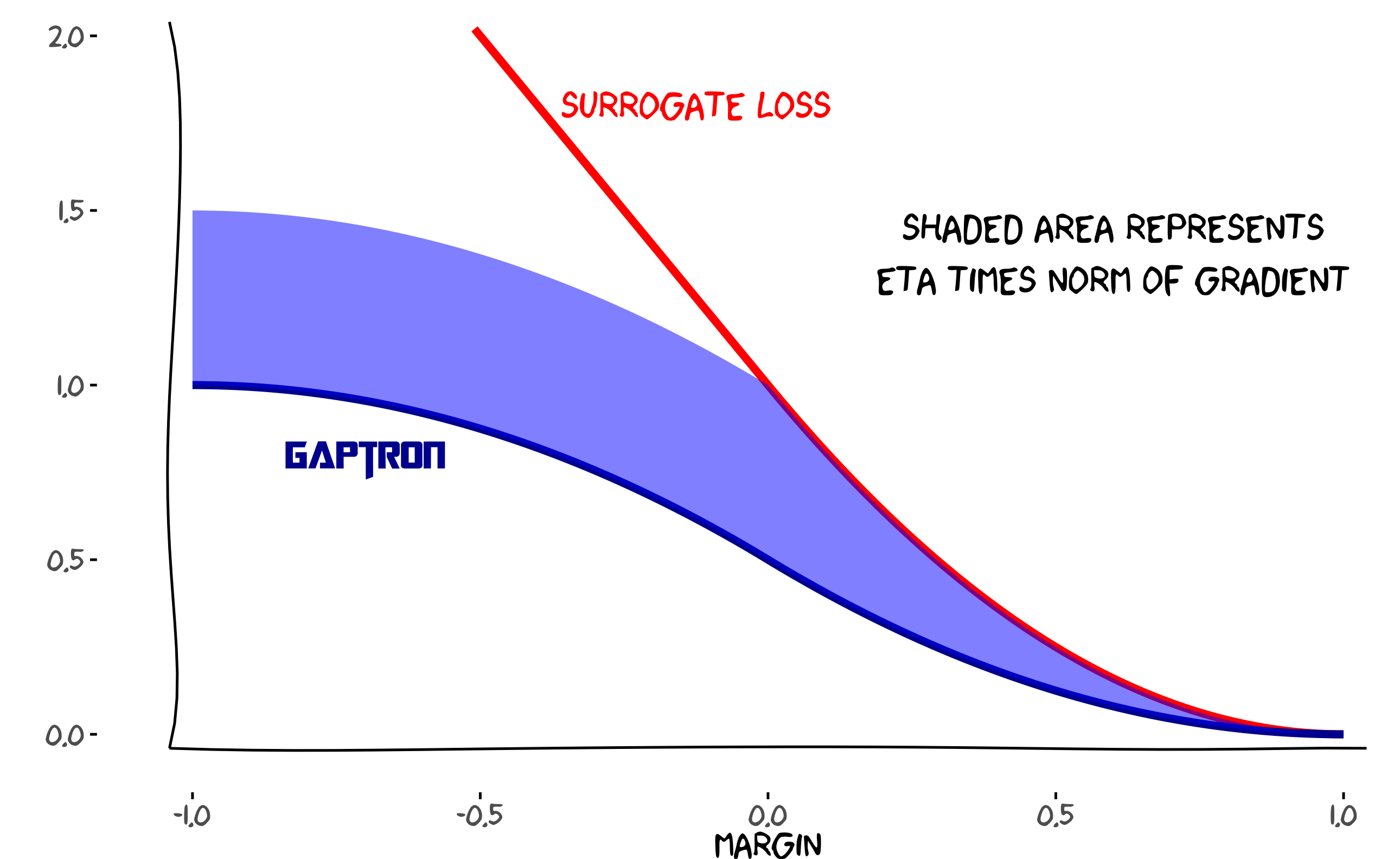


### Gaptron Key Idea



### Gaptron Analysis

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]] - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t) \\ &= \underbrace{\left( \sum_{t=1}^T \mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]] - \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) \right)}_{\text{Gaptron is random}} + \underbrace{\sum_{t=1}^T \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t)}_{\text{controlled by OGD}} \\ &\leq \frac{\|\mathbf{U}\|^2}{2\eta} + \left( \sum_{t=1}^T \mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]] + \frac{\eta}{2} \|\nabla_{\mathbf{W}_t} \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t)\|^2 - \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) \right) \end{aligned}$$



$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]] - \ell(\langle \mathbf{U}, \mathbf{x}_t \rangle, y_t) \\ &\leq \frac{\|\mathbf{U}\|^2}{2\eta} + \left( \sum_{t=1}^T \mathbb{E}[\mathbb{1}[y_t \neq \hat{y}_t]] + \frac{\eta}{2} \|\nabla_{\mathbf{W}_t} \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t)\|^2 - \ell(\langle \mathbf{W}_t, \mathbf{x}_t \rangle, y_t) \right) \\ &\leq 2KX^2\|\mathbf{U}\|^2 \end{aligned}$$

Smaller than surrogate loss!