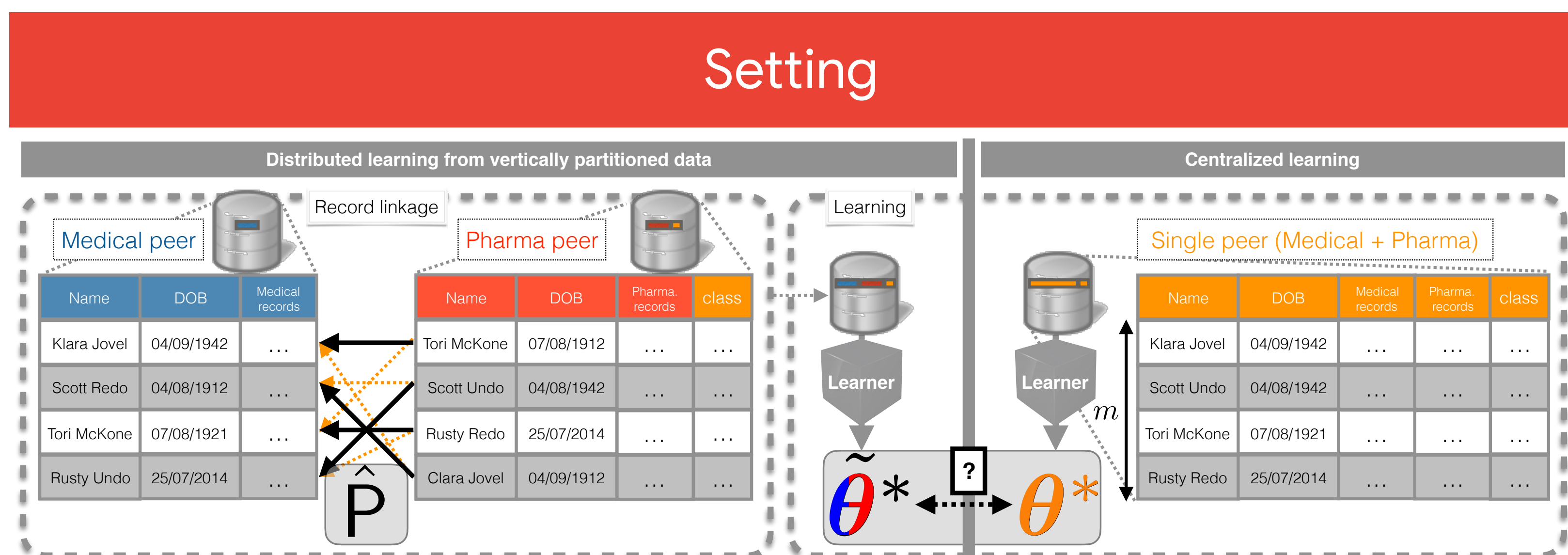


Richard Nock¹, Stephen Hardy², Wilko Henecka², Hamish Ivey-Law³,
Jakub Nabaglo³, Giorgio Patrini⁴, Guillaume Smith² & Brian Thorne²



- Our framework: vertical partition (VP):** features split among 2 peers, M and P (1+, e.g. P , holds label C)
- Record linkage (RL)** needed before batch supervised learning
- Baseline framework:** batch learning: labeled sample to learn classifier
- Example: **(Med+Pharma)** single **peer** holds all data

Problems, summarised: compare $\tilde{\theta}^*$ and θ^* ? Improve RL to get better $\tilde{\theta}^*$?

Summary

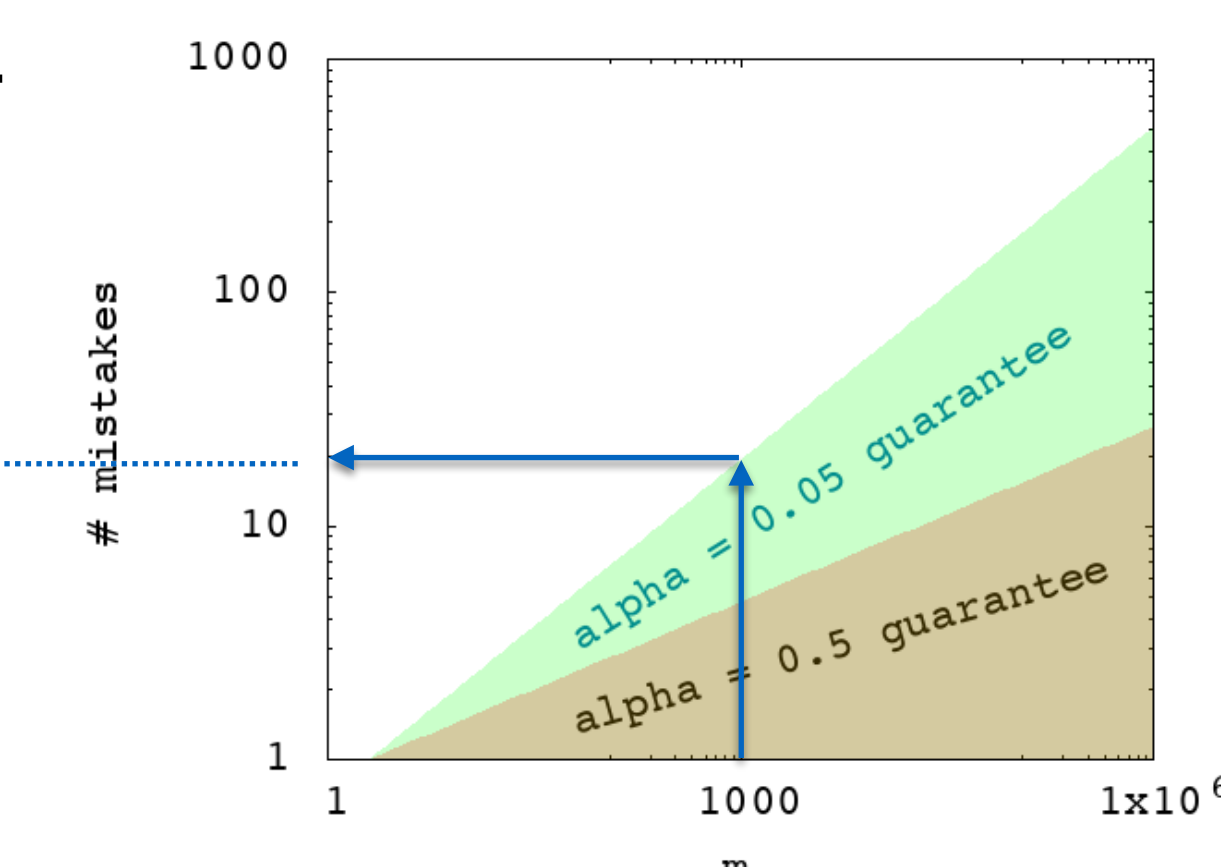
- Linear models: sufficient conditions to get on training:

$$\frac{\|\tilde{\theta}^* - \theta^*\|_2}{\|\theta^*\|_2} \leq O\left(\frac{1}{m^\alpha}\right)$$

- Regularisation is key (for broad set of losses)
- Optimisation of RL prior to ML, e.g. min. *between-classes* RL errors
- Large margins on $\theta^* \Rightarrow$ right class on $\tilde{\theta}^*$
- Results hold in the small data regime
- A small # of RL mistakes does not impact ML

“RL immunity for large margins”

RL does less than ~20 mistakes for 1000 ex. implies $\alpha = 0.05$.



- Exps include case where 1 peer does not hold class or a noisy estimate
- Simple approaches to complete / correct label prior to RL can \Rightarrow leverage
- Potential (estimated) **dials** to evaluate the value of RL prior to ML
 - inaccuracy between matched observations
 - error between classes for matched observations
- Margin immunity observed experimentally:**
 - strong advocacy for distributed / federated ML
 - may offer further cheap dials to evaluate the ML “potential” of datasets

Theory

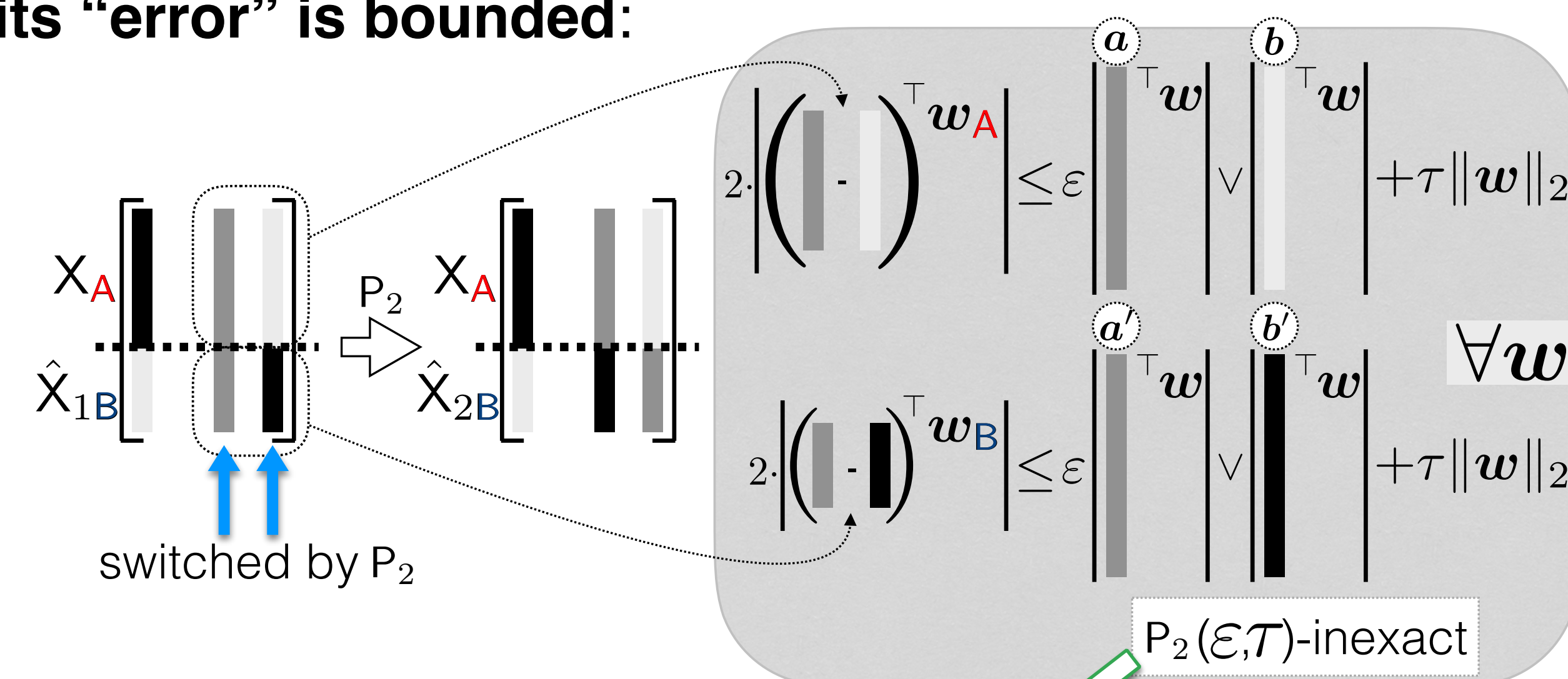
- Loss optimized** = Taylor loss + Ridge regularisation. For training sample $\hat{S} \doteq \{(\hat{x}_i, y_i), i \in [m]\}$, the loss is:

$$\ell_F(\hat{S}, \theta; \Gamma) \doteq \mathbb{E}_i[F(y_i \theta^\top \hat{x}_i)] + \theta^\top \Gamma \theta, \text{ with } F(z) \doteq a + bz + cz^2, \text{ and } a \in \mathbb{R}, b, c \in \mathbb{R}_*, \text{ and } \Gamma \text{ symmetric PSD}$$

- Record Linkage setting:** observation matrix X , split between two peers A and B. RL reconstructs \hat{X} but makes mistakes as (unknown) permutation matrix $\hat{P} \in \{0, 1\}^{m \times m}$

$$\hat{X} \doteq \begin{bmatrix} X_A \\ \hat{X}_B \doteq X_B \hat{P} \end{bmatrix} \xrightarrow{\text{(RL)}} \begin{bmatrix} X_A \\ X_B \end{bmatrix} \xrightarrow{\text{(unknown)}} X \doteq \begin{bmatrix} X_A \\ X_B \end{bmatrix}$$

Decomposition as elementary permutations: $\hat{P} = \prod_{t=1}^T P_t$
Each elementary permutation P_t switches parts between 2 observations. It is called (ϵ, τ) -inexact iff its “error” is bounded:



Key parameters for \hat{P} :
 $T = \text{size}$ and $\xi(\hat{P}) \doteq \min \left\{ \epsilon + \frac{\tau}{X_*} : \hat{P} \text{ is } (\epsilon, \tau)\text{-inexact} \right\}$
never ‘large’: $\xi(\cdot) \leq 2$ (all t)
max column-norm in X

Theorem: if ℓ_F ‘sufficiently’ regularised & m not too ‘small’,
 $\frac{\|\tilde{\theta}^* - \theta^*\|_2}{\|\theta^*\|_2} \leq \frac{\xi^{\frac{1}{4}}}{m} \cdot T^2 \cdot \left(\xi^{\frac{3}{4}} + \frac{\delta_{\hat{P}, l}}{\delta_\theta} \right)$
if $b = c$ and $m \geq 16$

where $\delta_\theta \doteq \|\theta^*\|_2 X_*$ and $\delta_{\hat{P}, l} \doteq \frac{\rho |b| L_\xi}{|c|}$ and $\rho \doteq \frac{T_+}{T} \in [0, 1]$

- Furthermore, $\forall \kappa > 0$, if m ‘large enough’:

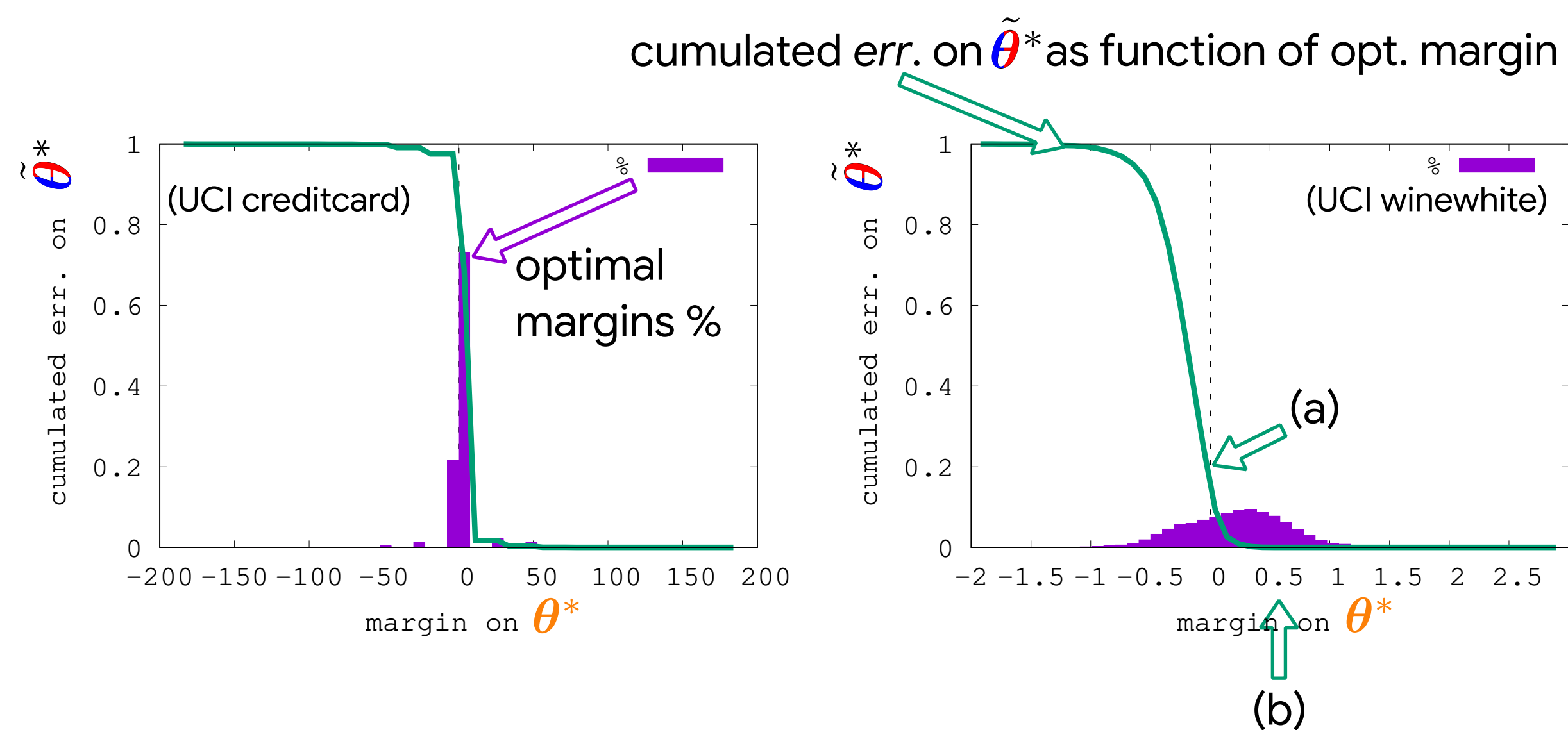
$$\frac{m}{\xi^{\frac{1}{4}} T^2} > \frac{\xi^{\frac{3}{4}} \delta_\theta + \delta_{\hat{P}, l}}{\kappa}$$

then $\forall (x, y), (y \theta^{*\top} x > \kappa) \Rightarrow (y \tilde{\theta}^{*\top} x > 0)$

(‘large’ optimal margin) \Rightarrow (right class despite RL mistakes)
 $(\tilde{\theta}^*$ is immune to record linkage at margin κ)

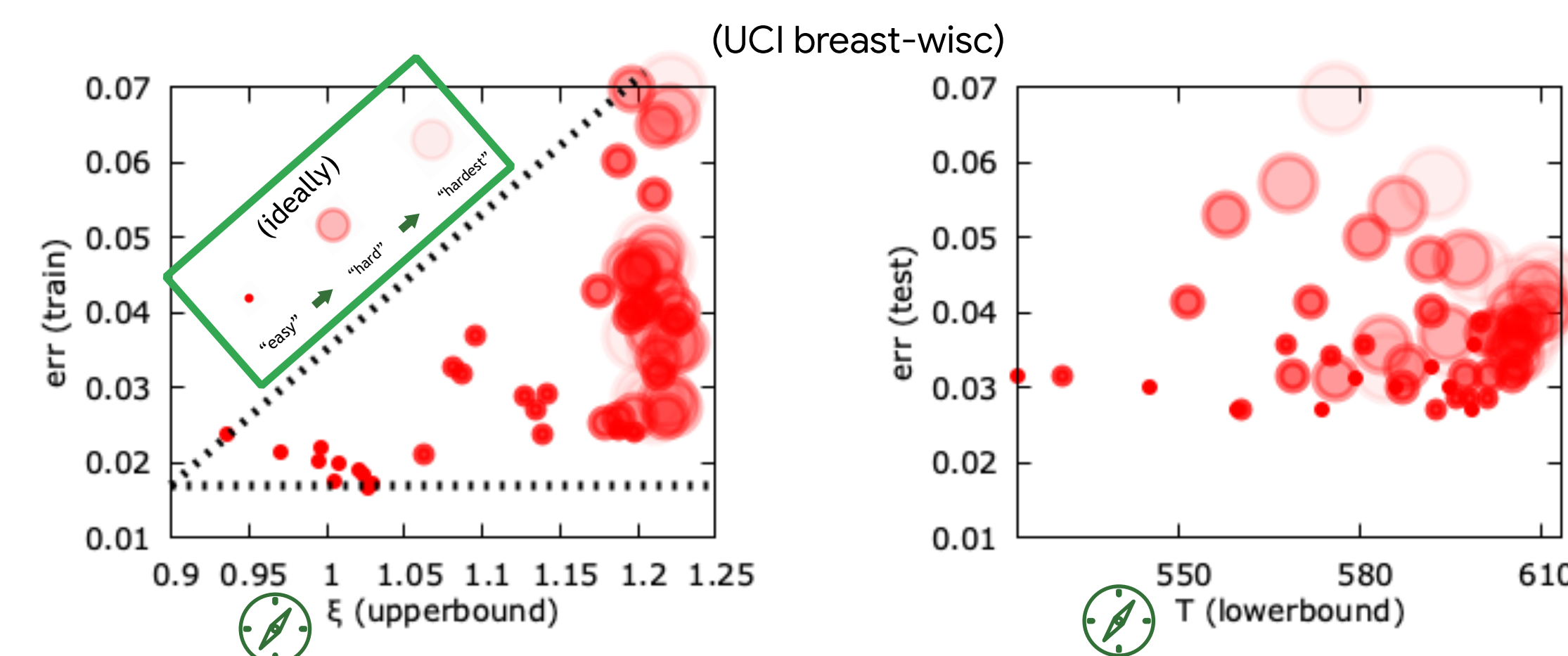
Experiments

- Experiments on RL immunity for large margin classification



(a): $\leq 20\%$ errors happen on examples with > 0 opt. margin
(b): almost **no error** on examples with optimal margin > 0.5
Justification for ML on VP setting: VP = more features, so opt. margins **increase**, beneficial for a **good** RL+ML pipeline

- Experiments on computable / estimable dials to ‘link’ the difficulty of RL to that of ML



Main parameters of \hat{P} = good dials for ML
Could be of use to select / improve RL techs

- Key parameters for \hat{P} :
- size T
 - “inaccuracy” ξ
 - “error” ρ

Theory

Experiments