

Beyond Empirical Risk Minimization

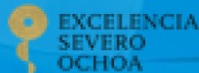
Santiago Mazuelas

smazuelas@bcamath.org

Basque Center for Applied Mathematics-BCAM

DATAI, University of Navarra, 21 of February 2024

(bcam)



www.bcamath.org
basque center for applied mathematics



Empirical Risk Minimization

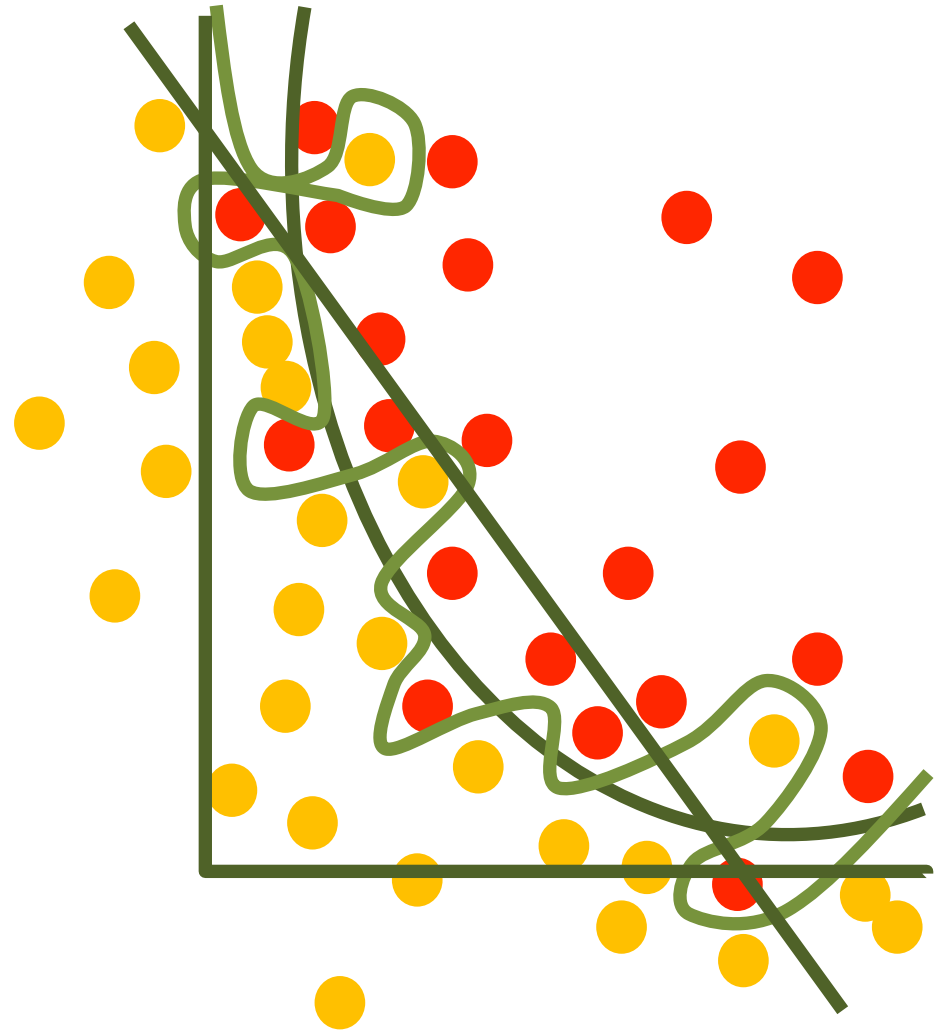
Consider multiple
“simple” hypothesis
and
choose one that fits well
the data



William of Ockham

Empirical Risk Minimization

Consider multiple
“simple” hypothesis
and
choose one that fits well
the data



Empirical Risk Minimization

Strongly rely on the specific training samples available

Not clear how to address corrupted samples or distribution shifts

The quantity optimized at learning is not very meaningful

Fitting training samples is not directly related with the prediction process



Robust Risk Minimization

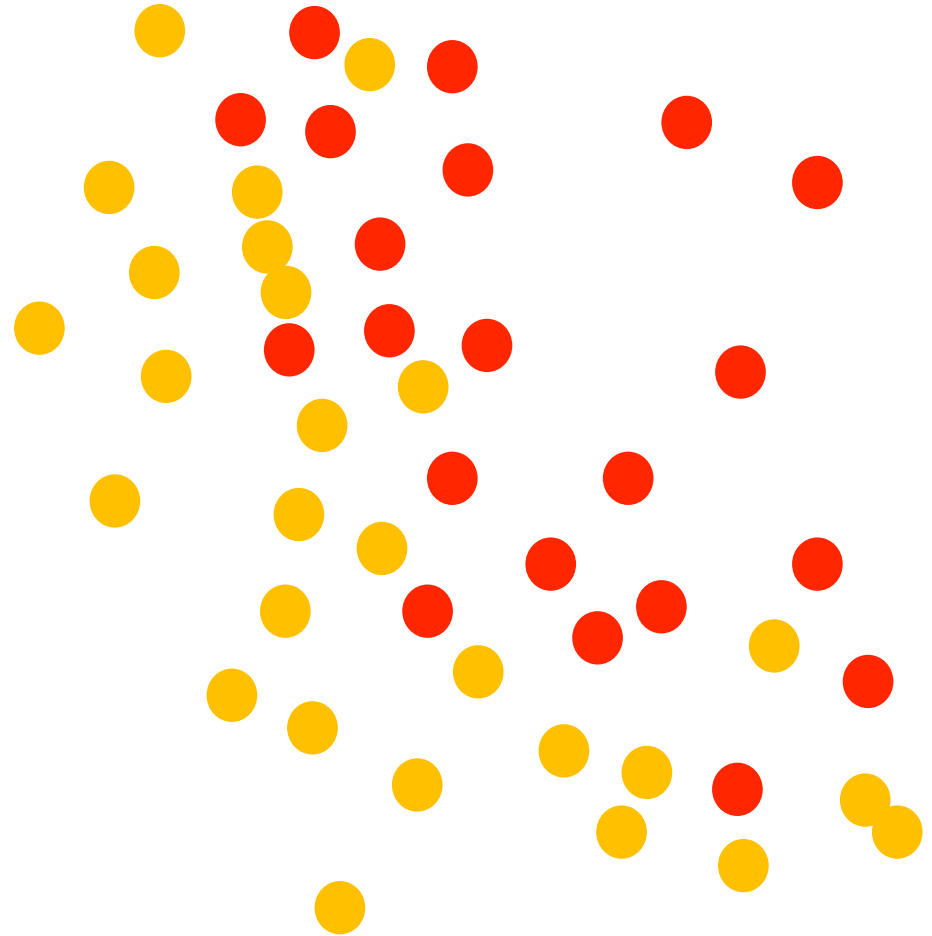
Consider multiple scenarios
consistent with the data
and
choose a rule with small error
in the worst scenario



John von Neumann

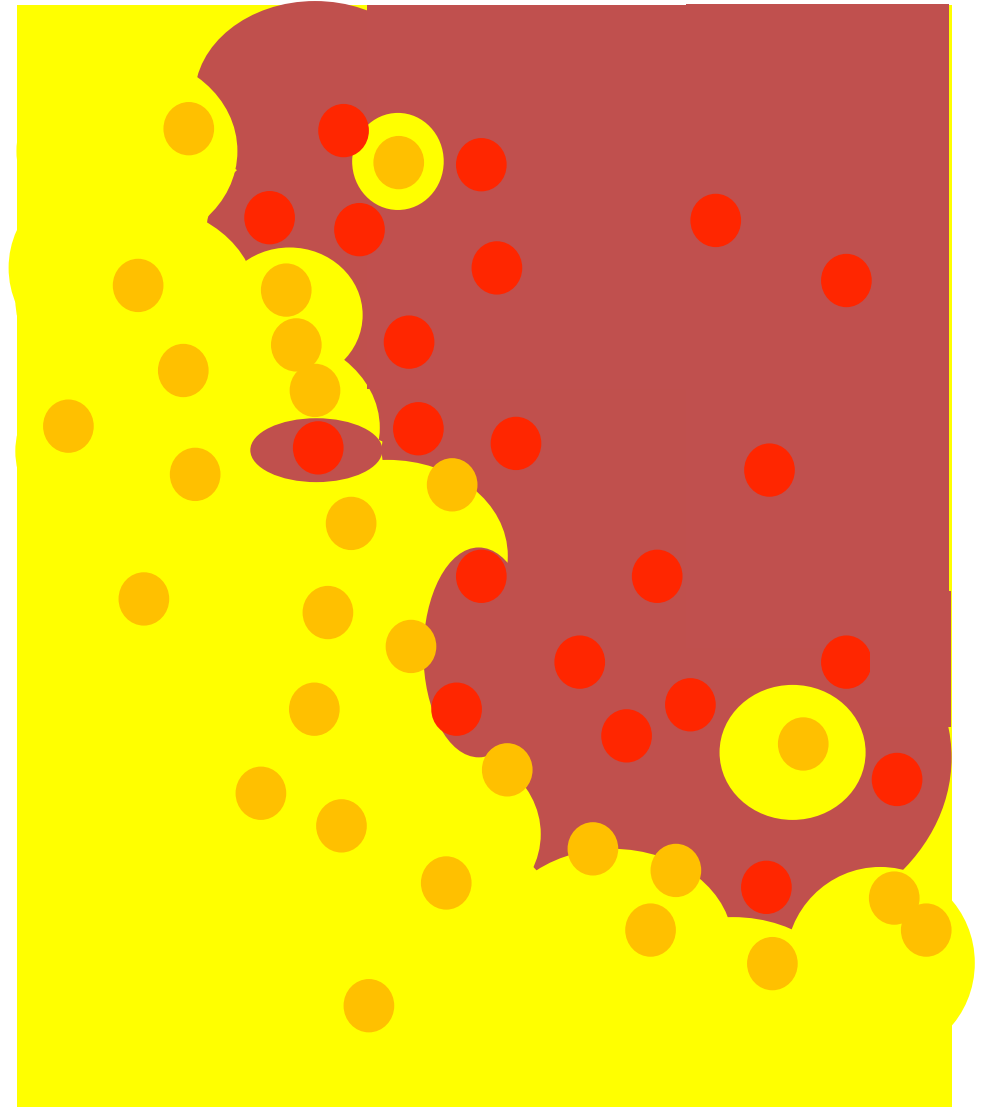
Robust Risk Minimization

Consider multiple scenarios
consistent with the data
and
choose a rule with small error
in the worst scenario



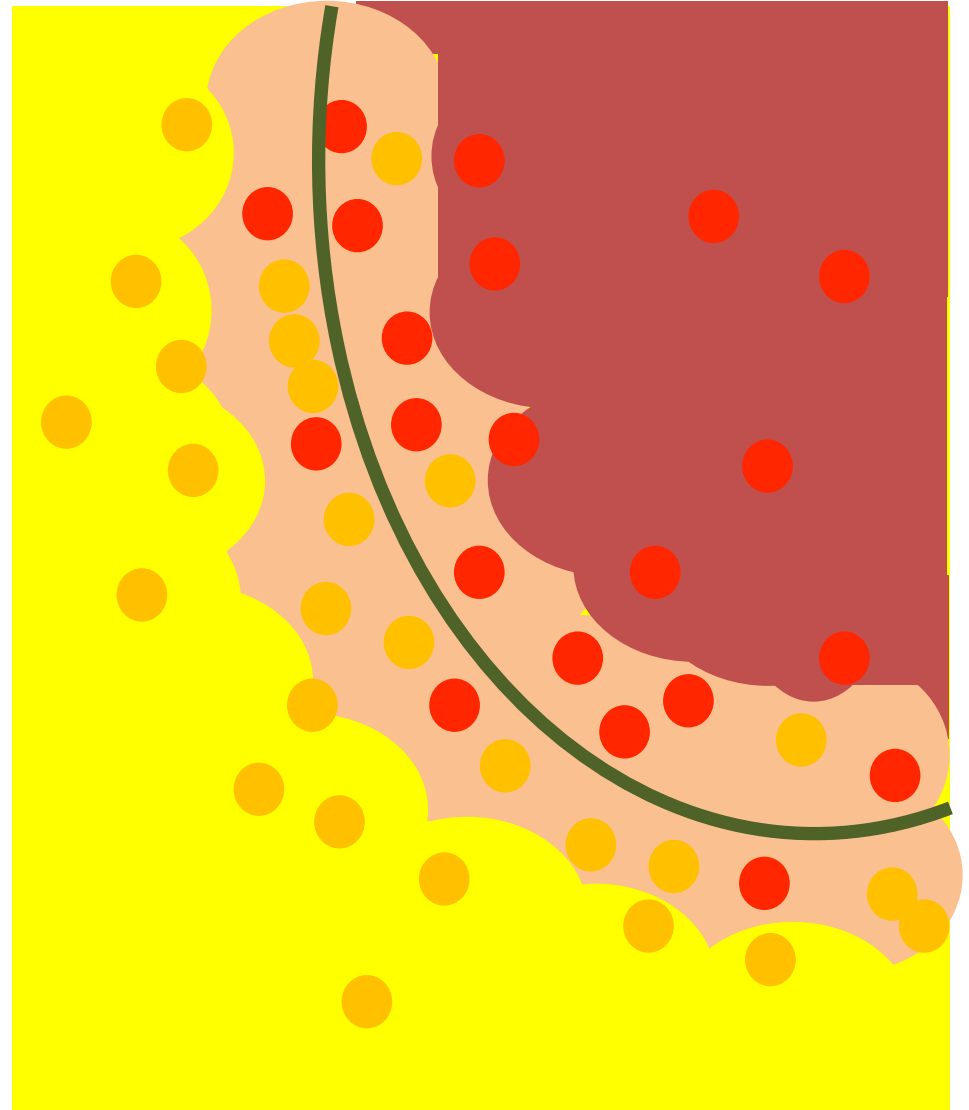
Robust Risk Minimization

Consider multiple scenarios
consistent with the data
and
choose a rule with small error
in the worst scenario



Robust Risk Minimization

Consider multiple scenarios
consistent with the data
and
choose a rule with small error
in the worst scenario



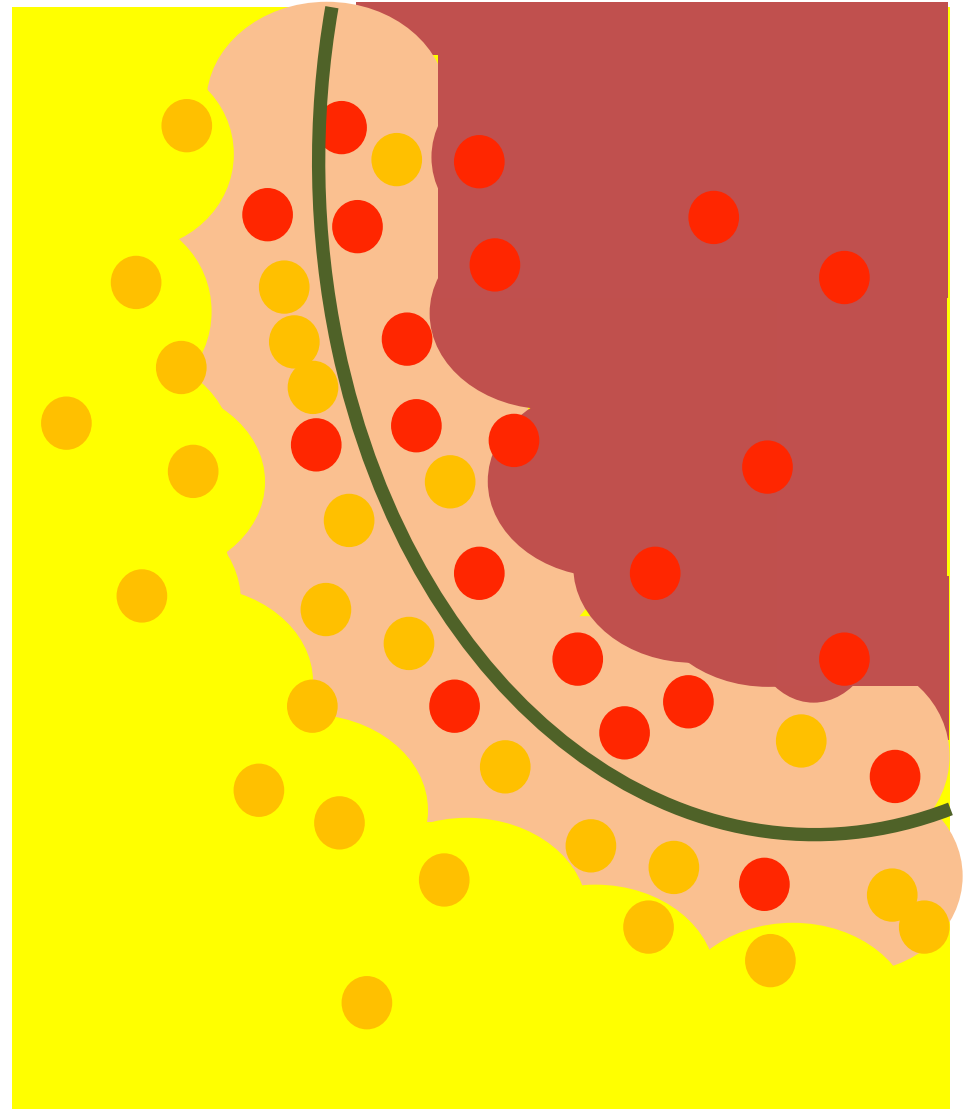
Robust Risk Minimization

Presence of corrupted samples or distribution shifts can be included on the considered scenarios

The quantity optimized at learning can provide an upper bound for the prediction error

Learning process analyzes the prediction performance of different rules

Can be too pessimistic with too broad feasible scenarios



Outline

- Introduction
- Generalized maximum entropy for supervised classification
- Minimax risk classifiers (MRCs)
- Performance guarantees 0-1 MRCs
- Multidimensional adaptation to concept drift

Maximum entropy principle

$$c : \text{Words} \rightarrow \{0, 1\}^*$$
$$\text{"cat"} \mapsto c(\text{"cat"}) = 1, 0, 1, 1$$

$$\text{Words} = \{w\} \quad \text{Codes} = \{c\} \quad \ell(c, w) = [c(w)]$$

$$w \sim p(w) \in \Delta(\text{Words})$$

$$\min_{c \in \text{Codes}} \mathbb{E}_p \ell(c, w) = \min_{c \in \text{Codes}} \ell(c, p) = -\mathbb{E}_p \log_2 p = H(p)$$

$$p \in \mathcal{U} \quad \text{Given by expectation constraints} \quad \rightarrow \max_{p \in \mathcal{U}} H(p)$$

[Jaynes, Phys. Rev., 1957], [P. Grünwald and A. P. Dawid, Annals of Statistics, 2004]

Supervised classification

Examples = $\{(x, y)\}$ x instance or attribute y label or class

Classification rules = $\{h : \mathcal{X} \rightarrow \Delta(\mathcal{Y})\} = \mathsf{T}(\mathcal{X}, \mathcal{Y})$

$$\ell_{0-1}(h, (x, y)) = \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{if } h(x) = y \end{cases} \quad \ell_{0-1}(h, (x, y)) = 1 - h(y|x)$$

$$\ell_{\log}(h, (x, y)) = -\log h(y|x)$$

ℓ -entropy for classification

$$H_{\ell}(p) = \min_{h \in \mathsf{T}(\mathcal{X}, \mathcal{Y})} \mathbb{E}_p \ell(h, (x, y)) = \min_{h \in \mathsf{T}(\mathcal{X}, \mathcal{Y})} \ell(h, p)$$

$$H_{0-1}(p) = 1 - \sum_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} p(x, y)$$

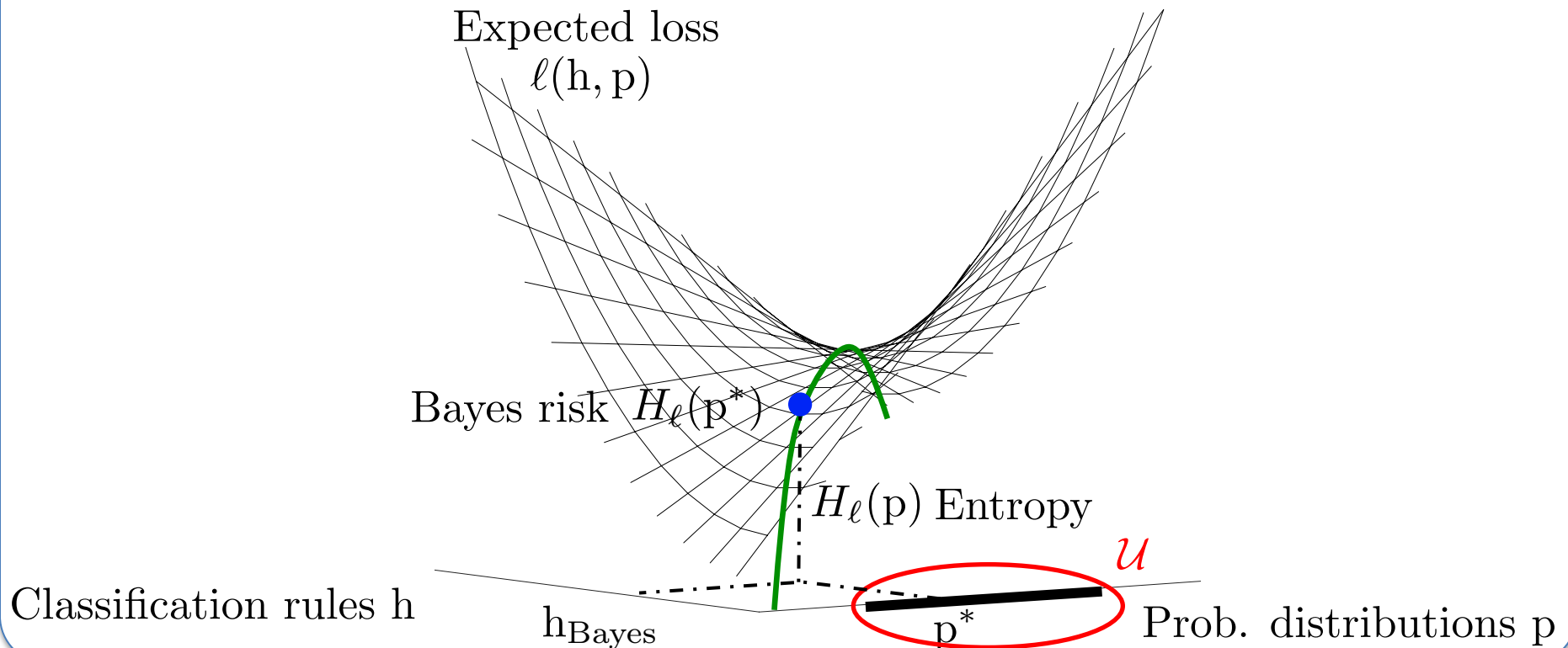
$$H_{\log}(p) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x)}{p(x, y)}$$

Maximum entropy and minimax rules

p^* underlying distribution of instance-label pairs

$$H_\ell(p^*) = \min_{h \in T(\mathcal{X}, \mathcal{Y})} \ell(h, p^*) \text{ Bayes risk}$$

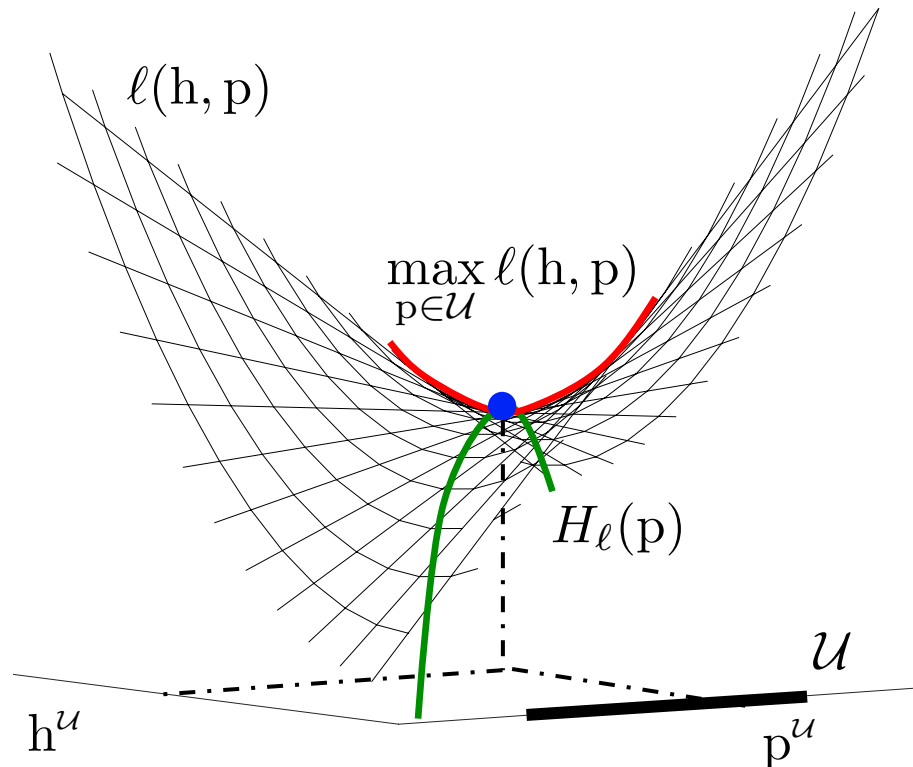
$$H_\ell(p^*) = \ell(h_{\text{Bayes}}, p^*), \quad h_{\text{Bayes}} \text{ Bayes classifier}$$



Maximum entropy and minimax rules

$p^u \in \arg \max_{p \in \mathcal{U}} H_\ell(p)$ Maximum entropy distribution

$h^u \in \arg \min_{h \in T(\mathcal{X}, \mathcal{Y})} \left(\max_{p \in \mathcal{U}} \ell(h, p) \right)$ Minimax risk classifier

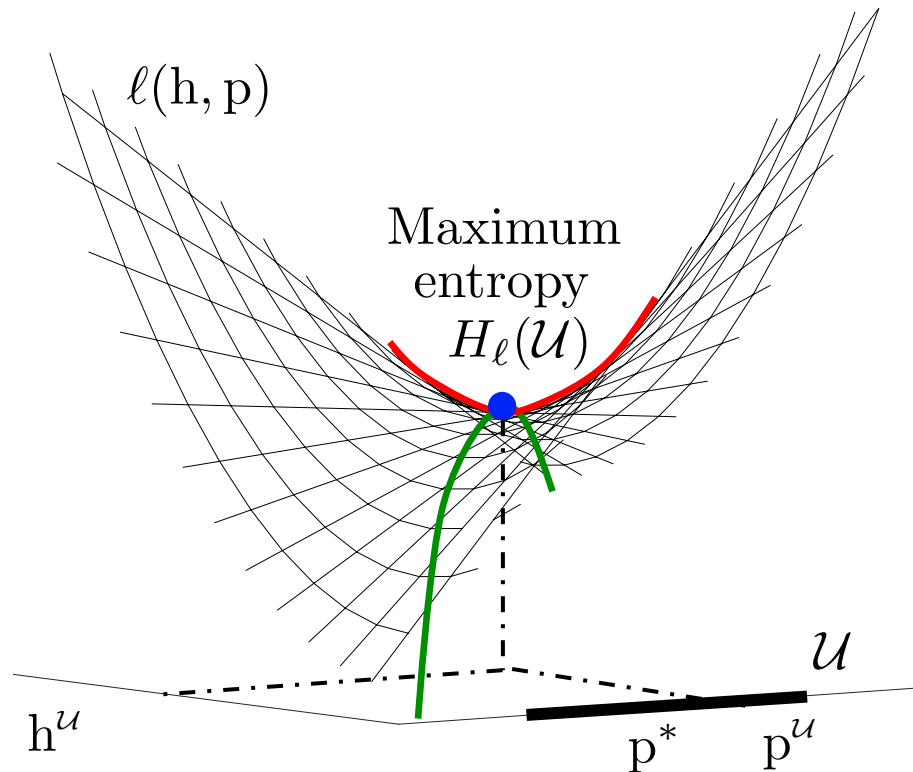


Maximum entropy and minimax rules

Theorem [S. Mazuelas, Y. Shen, A. Perez, IEEE Trans. Inf. Theory, 2022]

$$H_\ell(\mathcal{U}) = \max_{p \in \mathcal{U}} H_\ell(p) = \ell(h^u, p^u) = \min_{h \in T(\mathcal{X}, \mathcal{Y})} \left(\max_{p \in \mathcal{U}} \ell(h, p) \right)$$

$$h^u \in \arg \min_h \ell(h, p^u) \quad p^u \in \arg \max_p \ell(h^u, p)$$

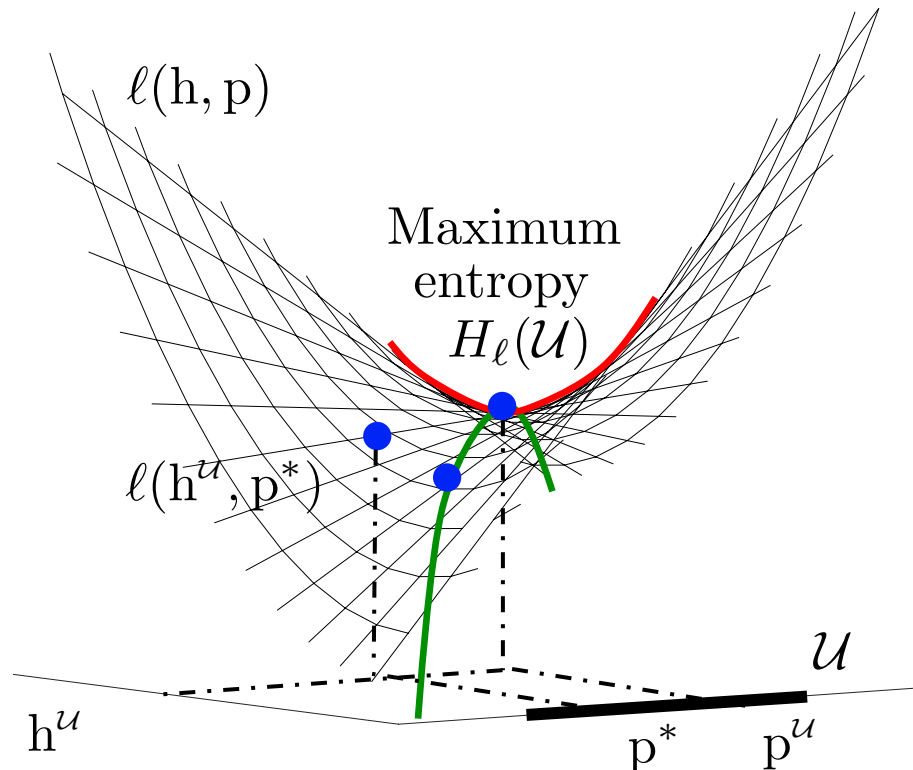


Maximum entropy and minimax rules

Theorem [S. Mazuelas, Y. Shen, A. Perez, IEEE Trans. Inf. Theory, 2022]

$$H_\ell(\mathcal{U}) = \max_{p \in \mathcal{U}} H_\ell(p) = \ell(h^u, p^u) = \min_{h \in T(\mathcal{X}, \mathcal{Y})} \left(\max_{p \in \mathcal{U}} \ell(h, p) \right)$$

$$p^* \in \mathcal{U} \Rightarrow H_\ell(p^*) \leq \ell(h^u, p^*) \leq H_\ell(\mathcal{U})$$



Outline

- Introduction
- Generalized maximum entropy for supervised classification
- **Minimax risk classifiers (MRCs)**

Minimax Risk Classifiers (MRCs)

$$\mathcal{P}_{\text{MRC}} : \min_{h \in T(\mathcal{X}, \mathcal{Y})} \max_{p \in \mathcal{U}} \ell(h, p)$$

$$\mathcal{U} = \left\{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p\{\Phi\} - \boldsymbol{\tau}| \preceq \boldsymbol{\lambda} \right\}$$

$\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$ feature mapping

$$\Phi(x, y) = \mathbf{e}_y \otimes \boldsymbol{\Psi}(x), \quad \boldsymbol{\Psi}(x) = [\psi_{v_1}(x), \psi_{v_2}(x), \dots, \psi_{v_D}(x)]^T$$

e.g., random features $\mathbb{E}_{p(v)}\{\psi_v(x)\psi_v(x')\} = k(x, x')$

$$\boldsymbol{\tau}_n = \frac{1}{n} \sum_{i=1}^n \Phi(x_i, y_i)$$

$\boldsymbol{\lambda}$ confidence vector at level $1 - \delta \Rightarrow p^* \in \mathcal{U}$ w.p. $\geq 1 - \delta$

MRCs Learning. Representer Theorem

Theorem (0-1 loss) [S. Mazuelas, A. Zanoni, A. Perez, NeurIPS 2021]

Let $\mu^* \in \mathbb{R}^m$ be a solution of

$$\min_{\mu \in \mathbb{R}^m} 1 - \tau^T \mu + \varphi(\mu) + \lambda^T |\mu|$$

where $\varphi(\mu) = \max_{x \in \mathcal{X}, \mathcal{C} \subseteq \mathcal{Y}} \frac{\sum_{y \in \mathcal{C}} \Phi(x, y)^T \mu - 1}{|\mathcal{C}|}$.

If $h^u \in T(\mathcal{X}, \mathcal{Y})$ satisfies

$$h^u(y|x) \geq \Phi(x, y)^T \mu^* - \varphi(\mu^*), \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}$$

Then

$$h^u \in \arg \min_{h \in T(\mathcal{X}, \mathcal{Y})} \max_{p \in \mathcal{U}} \ell(h, p)$$

$$H_{0-1}(\mathcal{U}) = \bar{R}(\mathcal{U}) = 1 - \tau^T \mu^* + \varphi(\mu^*) + \lambda^T |\mu^*|$$

Outline

- Introduction
- Generalized maximum entropy for supervised classification
- Minimax risk classifiers (MRCs)
- Performance guarantees 0-1 MRCs

MRCs Generalization Bounds

$$R_{\Phi} = \min_{\mu} 1 - \mathbb{E}_{p^*} \Phi(x, y)^T \mu + \varphi(\mu) = 1 - \mathbb{E}_{p^*} \Phi(x, y)^T \mu_{\infty} + \varphi(\mu_{\infty})$$

The MRC given by μ_{∞} is the optimal minimax rule for Φ :

$$R_{\Phi} \leq \bar{R}(U) \quad \forall U = \{p \in \Delta(\mathcal{X}, \mathcal{Y}) : |\mathbb{E}_p\{\Phi\} - \tau| \leq \lambda\} \text{ s.t. } p^* \in U$$

Theorem [S. Mazuelas, M. Romero, P. Grünwald, JMLR 2023]

$$R(h^U) \leq \bar{R}(U) + (|\mathbb{E}_{p^*} \Phi - \tau| - \lambda)^T |\mu^*|$$

$$R(h^U) \leq R_{\Phi} + |\mathbb{E}_{p^*} \Phi - \tau|^T |\mu_{\infty} - \mu^*| + \lambda^T (|\mu_{\infty}| - |\mu^*|)$$

In particular, if $\mathbb{P}\{|\mathbb{E}_{p^*} \Phi - \tau| \leq \lambda\} \geq 1 - \delta$

$$R(h^U) \leq \bar{R}(U)$$

$$R(h^U) \leq R_{\Phi} + 2\lambda^T |\mu_{\infty}|$$

w. p. at least $1 - \delta$.

MRCs Universal Consistency

Theorem [S. Mazuelas, M. Romero, P. Grünwald, JMLR 2023]

Let Φ_n be a feature mapping given by D_n random features corresponding with a characteristic kernel.

If $\lim_{n \rightarrow \infty} D_n = \infty$, for any underlying distribution p^* we have that

$$\lim_{n \rightarrow \infty} R_{\Phi_n} = R_{\text{Bayes}}$$

with probability 1.

If h_n is the MRC given by n samples and

$$\lim_{n \rightarrow \infty} \lambda_n = 0 \quad \lim_{n \rightarrow \infty} \lambda_n \sqrt{n / \log n} = \infty$$

for any underlying distribution p^* we have that

$$\lim_{n \rightarrow \infty} R(h_n) = R_{\text{Bayes}}$$

with probability 1.

Outline

- Introduction
- Generalized maximum entropy for supervised classification
- Minimax risk classifiers (MRCs)
- Performance guarantees 0-1 MRCs
- **Multidimensional adaptation to concept drift**

Supervised classification under concept drift

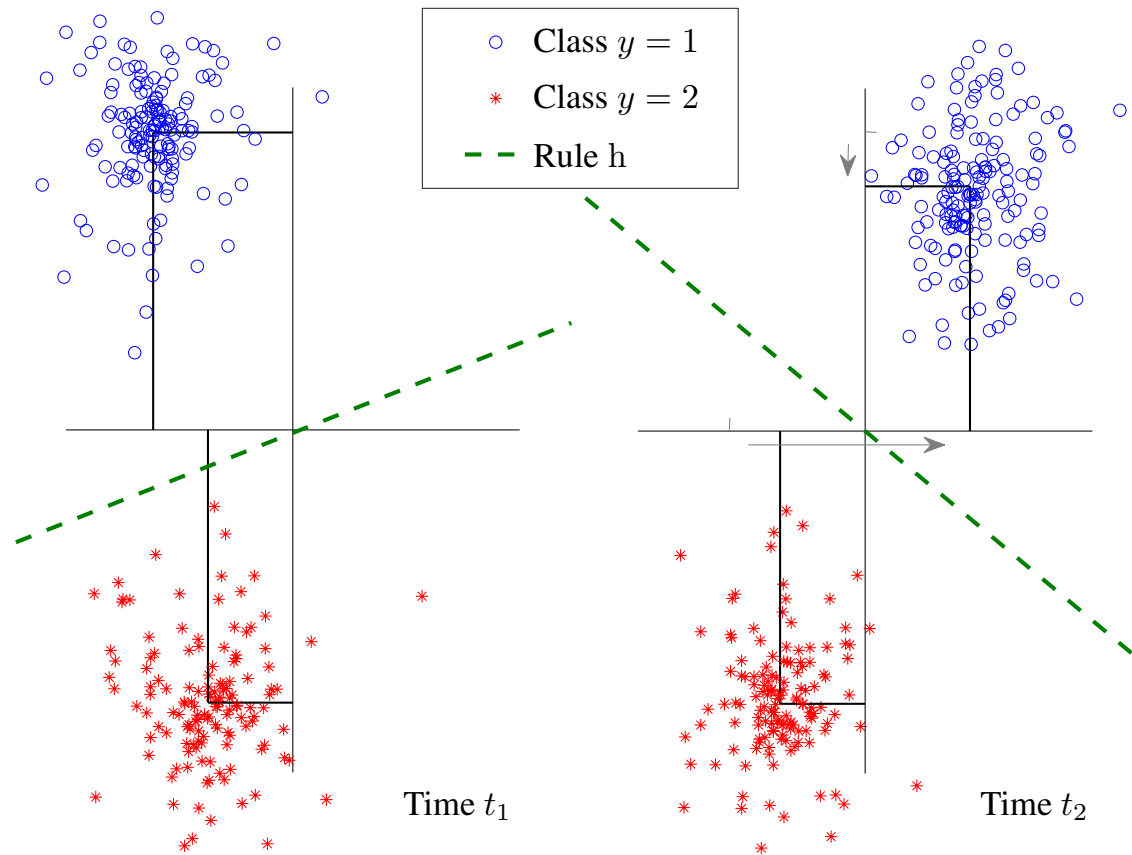
$$\begin{array}{cccccc} p_1 & & p_2 & & \dots & & p_t & & & & p_{t+1} & & \dots \\ \downarrow & & \downarrow & & & & \downarrow & & & & \downarrow & & \\ (x_1, y_1), & (x_2, y_2), & \dots, & (x_t, y_t), & (x_{t+1}, y_{t+1}), & \dots \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \\ h_1 & \rightarrow & h_2 & & \dots & & h_t & \rightarrow & h_{t+1} & & \dots \end{array}$$

For instance,

$$h_{t+1} = h_t - k \nabla \ell(h_t, (x_t, y_t))$$

$k \in \mathbb{R}$ accounts for a global rate of change

Supervised classification under concept drift



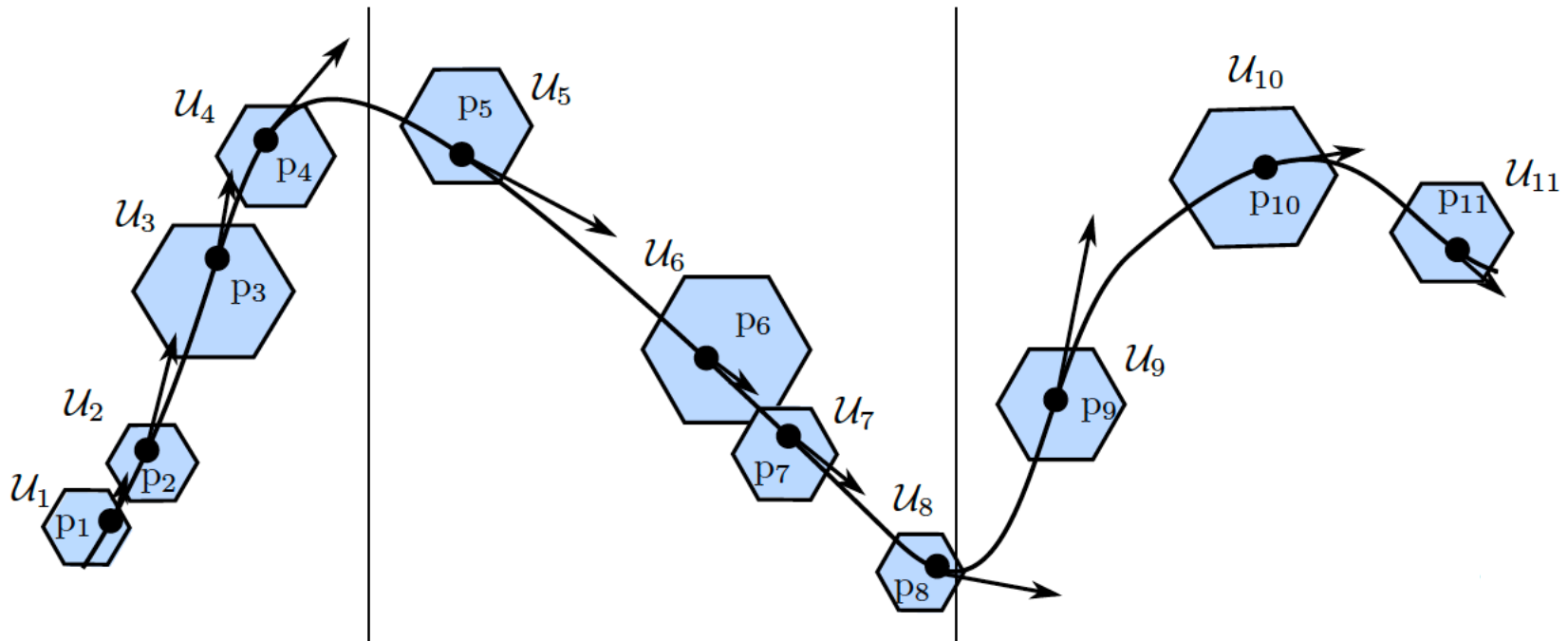
$$h_{t+1} = h_t - \boxed{k} \nabla \ell(h_t, (x_t, y_t))$$

$k \in \mathbb{R}$ accounts for a global rate of change

MRCs for concept drift adaptation

$$\mathcal{U}_t = \left\{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p\{\Phi\} - \hat{\tau}_t| \preceq \lambda_t \right\}$$

$$\hat{\tau}_t \approx \hat{\tau}_{t-1} \in \mathbb{R}^m$$



Multidimensional Adaptation

[V. Alvarez, S. Mazuelas, J. A. Lozano, ICML 2022]

[V. Alvarez, S. Mazuelas, J. A. Lozano, NeurIPS 2023]

Tracking underlying distribution: obtain $\hat{\tau}_t, \lambda_t$

Dynamical system: model the evolution of each component of τ_t

$$\eta_{t,i} = \mathbf{H}_t \eta_{t-1,i} + \mathbf{w}_{t,i} \quad \text{with } \eta_{t,i} = [\tau_{t,i}, \tau'_{t,i}, \tau''_{t,i}, \dots, \tau_{t,i}^{(k)}]^T$$
$$\Phi_i(x_t, y_t) = \tau_{t,i} + v_{t,i}$$

Unbiased linear estimator of $\eta_{t,i}$ with minimum MSE

$$\hat{\eta}_{t,i} = \mathbf{H}_t \hat{\eta}_{t-1,i} - \mathbf{k}_{t,i} (\hat{\tau}_{t-1,i} - \Phi_i(x_{t-1}, y_{t-1}))$$

$\mathbf{k}_{t,1}, \mathbf{k}_{t,2}, \dots, \mathbf{k}_{t,m}$ account for multidimensional time changes

Performance guarantees

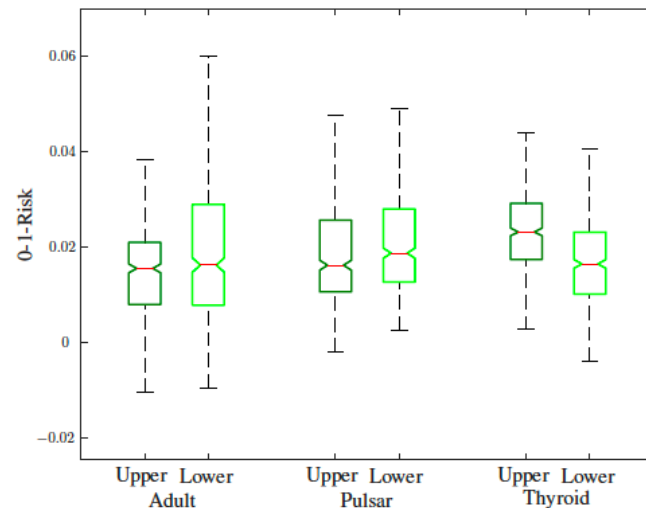
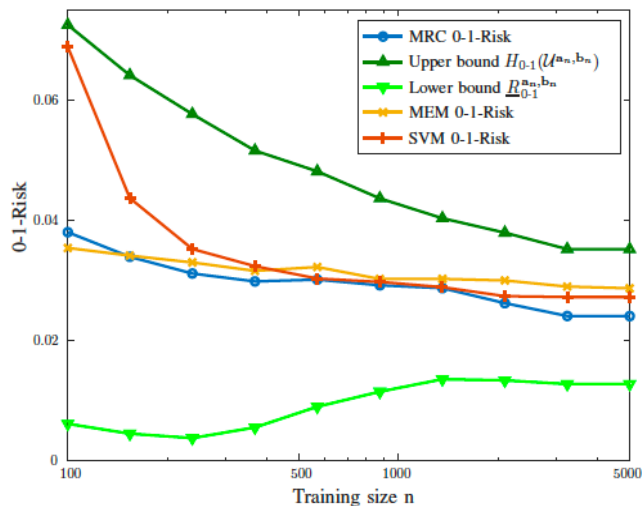
Error probability

$$R(h_t) \leq R(\mathcal{U}_t)$$

Accumulated mistakes

$$\sum_{t=1}^T \mathbb{I}\{\hat{y}_t \neq y_t\} \leq \sum_{t=1}^T R(\mathcal{U}_t) + \sqrt{2T \log \frac{1}{\delta}}$$

Some numerical results



<https://github.com/MachineLearningBCAM/MRCpy> <https://machinelearningbcam.github.io/MRCpy/>

Algorithm	Weather	Elec2	German	Chess	Usenet1	Email	Poker
DWM	30.0	36.3	41.2	35.2	36.3	39.5	22.0
Projectron	30.7	36.8	40.0	35.1	49.1	48.2	22.6
Unid. AMRC	31.3	40.1	30.3	38.3	46.3	48.4	39.4
AMRC	32.3	35.8	30.3	27.7	35.7	43.7	21.9
Det. AMRC	30.0	33.9	30.0	33.4	32.0	33.9	21.9

Conclusions

- Learning methodology that **relies on expectation estimates** and not on specific training samples
- MRCs can minimize the worst-case expected **0-1 loss** over general classification rules, and provide **performance guarantees** during learning
- Described MRCs' **generalization bounds** and universal consistency
- The proposed methods can lead to algorithms that **seamlessly address distribution shifts**

References

- S. Mazuelas, A. Zanoni, and A. Perez, “Minimax classification with 0-1 loss and performance guarantees,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 302–312.
- S. Mazuelas, Y. Shen, and A. Perez, “Generalized maximum entropy for supervised classification,” *IEEE Transactions on Information Theory*, vol. 68, no. 4, pp. 2530–2550, Apr. 2022.
- V. Alvarez, S. Mazuelas, and J. A. Lozano, “Minimax classification under concept drift with multidimensional adaptation and performance guarantees,” *International Conference on Machine Learning (ICML) 2022*.
- S. Mazuelas, M. Romero, and P. Grünwald, “Minimax classification with 0-1 loss,” *Journal of Machine Learning Research (JMLR)*, vol. 24, no. 208, pp. 1-48, 2023.
- V. Alvarez, S. Mazuelas, and J. A. Lozano, “Minimax Forward and Backward Learning of Evolving Tasks with Performance Guarantees,” *Advances in Neural Information Processing Systems (NeurIPS)*. 2023.
- J. Segovia, S. Mazuelas, and A. Liu, “Double-weighting for covariate shift adaptation,” *International Conference on Machine Learning (ICML) 2023*.