# Learning-based decision-making via sample compression: theoretical results and algorithms

# speaker: Simone Garatti

(Politecnico di Milano, Italy – email: simone.garatti@polimi.it)





## Many thanks to all collaborators!





#### Marco C. Campi

# Many thanks to all collaborators!





#### Marco C. Campi



Algo Carè



Federico Ramponi



Kostas Margellos



Alessandro Falsone



Maria Prandini



Dario Paccagnan

# Many thanks to all collaborators!





#### Marco C. Campi



Algo Carè



Federico Ramponi



Kostas Margellos



Alessandro Falsone



Maria Prandini



## **Learning-based decision-making**



# $\delta$ = uncertain element $\implies$ exercise caution



## **Learning-based decision-making**



# $\delta$ = uncertain element $\implies$ exercise caution



## **Learning-based decision-making**







## **Example: classification**



#### data

 $\delta_i = (x_i, y_i)$  $x_i \in \mathbb{R}^d$  $y_i \in \{\mathsf{red}, \mathsf{blue}\}$ 



### **Example: classification**



#### data

 $\delta_i = (x_i, y_i)$  $x_i \in \mathbb{R}^d$  $y_i \in \{\mathsf{red}, \mathsf{blue}\}$ 



 $\mathcal{H}$  = SVM classifier

$$\min_{\substack{w \in \Phi, b \in \mathbb{R} \\ \varepsilon_i \ge 0, i=1,\dots, N}} \|w\|^2 + \rho \sum_{i=1}^N \xi_i$$
  
s.t.  $1 - y_i(\langle w, \phi(x_i) \rangle - b) \le \xi_i, \quad i = 1,\dots, N$ 

# **Example: scenario robust optimization**







# **Example: scenario robust optimization**





# **Example:** scenario robust optimization (cont'd)



rate of return of k financial assets in day i  $\delta_i = (R_i^1, R_i^2, \dots, R_i^k)$ 

money split over the k financial assets

 $\theta = (\$^1, \$^2, \dots, \$^k, L)$ 

$$\begin{array}{ccc} \min_{\substack{\$^{i},L\\\$^{1}+\cdots+\$^{k}=1}} & L\\ \text{s.t.} & -(\$^{1}R_{i}^{1}+\$^{2}R_{i}^{2}+\cdots+\$^{k}R_{i}^{k}) \leq L\\ & i=1,\ldots,N \end{array}$$

# **Example: scenario optimization with relaxation**





# **Example: scenario optimization with relaxation**







#### Which is the learning-based scheme for the problem at hand?

Difficult to say a-priori without incurring in over-conservatism ... a blend of approximate knowledge and heuristics, often in various attempts (hyperparameters tuning)

No limits in exploration, but some guidance is needed...





 $\mathsf{R}(\mathcal{H}) = \mathbb{P}_{\delta} \{ \mathcal{H} \text{ is inappropriate for a new} \delta \}$   $\downarrow$ interaction between decision and environment

For instance, in scenario optimization:

 $\mathbb{P}_{\delta}$  { a new constraint is violated by  $\theta^*$  }







For instance, in scenario optimization:

 $\mathbb{P}_{\delta}$  { a new constraint is violated by  $\theta^*$  }



**Issue:**  $\mathbb{P}$  is not available...

is it possible to assess R(H) from data ?



- holding out some data for testing rather than designing... waste of information, questionable!
- scenarios (data) are often limited resources (collecting data can be time-consuming or burdensome, involving a monetary cost)
- in the context of this talk, testing is not necessary...

... data can well play a double role!



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



#### **Preference**

$$\kappa(\delta_1, \dots, \delta_N) \subseteq S \subseteq (\delta_1, \dots, \delta_N)$$
$$\bigvee_{\kappa(S) = \kappa(\delta_1, \dots, \delta_N)}$$



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



#### **Preference**

$$\kappa(\delta_1, \dots, \delta_N) \subseteq S \subseteq (\delta_1, \dots, \delta_N)$$
$$\bigvee_{\kappa(S) = \kappa(\delta_1, \dots, \delta_N)}$$



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



#### **Preference**

$$\kappa(\delta_1, \dots, \delta_N) \subseteq S \subseteq (\delta_1, \dots, \delta_N)$$
$$\bigvee_{\kappa(S) = \kappa(\delta_1, \dots, \delta_N)}$$



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



#### **Coherence**

a new scenario for which  ${\mathcal H}$  is inappropriate is added  $\bigcup_V$ 

the compression must change



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



#### **Coherence**

a new scenario for which  ${\mathcal H}$  is inappropriate is added  $\bigcup_V$ 

the compression must change



$$\kappa(\delta_1,\ldots,\delta_N)=\delta_{i_1},\ldots,\delta_{i_k}$$

map extracting a data subsample from the dataset



## **Coherence**

a new scenario for which  ${\mathcal H}$  is inappropriate is added  $\bigcup_V$ 

the compression must change

#### The main result in a nutshell



**Risk:** 
$$\mathsf{R}(\mathcal{H}) = \mathsf{R}(\mathcal{H}(\delta_1, \dots, \delta_N))$$



Risk:
$$\mathsf{R}(\mathcal{H}) = \mathsf{R}(\mathcal{H}(\delta_1, \dots, \delta_N))$$
randomComplexity: $\pi = |\kappa(\delta_1, \dots, \delta_N)|$ variables $\uparrow$ size of compressed set



Risk:
$$\mathsf{R}(\mathcal{H}) = \mathsf{R}(\mathcal{H}(\delta_1, \dots, \delta_N))$$
randomComplexity: $\pi = |\kappa(\delta_1, \dots, \delta_N)|$ variables

Under preference and coherence, the joint distribution of risk and complexity is concentrated around/below  $R(H) = \pi/N$ 





Risk:
$$\mathsf{R}(\mathcal{H}) = \mathsf{R}(\mathcal{H}(\delta_1, \dots, \delta_N))$$
randomComplexity: $\pi = |\kappa(\delta_1, \dots, \delta_N)|$ variables

Under preference and coherence, the joint distribution of risk and complexity is concentrated around/below  $R(H) = \pi/N$ 



$$\mathsf{R}(\mathcal{H})$$
 can be accurately estimated from  $\pi$ 



Risk:
$$\mathsf{R}(\mathcal{H}) = \mathsf{R}(\mathcal{H}(\delta_1, \dots, \delta_N))$$
randomComplexity: $\pi = |\kappa(\delta_1, \dots, \delta_N)|$ variables

Under preference and coherence, the joint distribution of risk and complexity is concentrated around/below  $R(H) = \pi/N$ 





**Theorem** (with M. Campi) Assume preference and coherence Choose  $\beta \in (0,1)$  (confidence parameter) Let  $\epsilon^{U}(k)$  be the unique roots in (0,1) of polynomials  $\triangleright \binom{N}{k} (1-\epsilon)^{N-k} - \frac{\beta}{2N} \sum_{i=1}^{N-1} \binom{m}{k} (1-\epsilon)^{m-k}$ Then, irrespective of  $\mathbb{P}$  (distribution-free),  $\mathbb{P}^{N}\left\{\delta_{1},\ldots,\delta_{N}:\ \mathsf{R}(\mathcal{H})\leq\epsilon^{U}(\pi)\right\}\geq1-\beta$ 















# $\mathsf{R}(\mathcal{H}) \leq \epsilon^U(\pi)$ is true with confidence $1 - \beta$





# $\mathsf{R}(\mathcal{H}) \leq \epsilon^U(\pi)$ is true with confidence $1 - \beta$





# $\mathsf{R}(\mathcal{H}) \leq \epsilon^{U}(\pi)$ is true with confidence $1 - \beta$





#### Theorem (with M. Campi)

Assume preference and coherence + additional assumptions

Choose  $\beta \in (0,1)$  (confidence parameter)

Let  $\epsilon_L(k), \epsilon^U(k)$  be the unique roots in (0,1) of polynomials

$$\triangleright \quad \binom{N}{k} (1-\epsilon)^{N-k} - \frac{\beta}{2N} \sum_{m=k}^{N-1} \binom{m}{k} (1-\epsilon)^{m-k}$$
$$\triangleright \quad \binom{N}{k} (1-\epsilon)^{N-k} - \frac{\beta}{2N} \sum_{m=N+1}^{2N} \binom{m}{k} (1-\epsilon)^{m-k}$$

Then, irrespective of  $\mathbb{P}$ ,

 $\mathbb{P}^{N}\left\{\delta_{1},\ldots,\delta_{N}:\ \epsilon_{L}(\pi)\leq\mathsf{R}(\mathcal{H})\leq\epsilon^{U}(\pi)\right\}\geq1-\beta$ 



# $\epsilon_L(\pi) \leq \mathsf{R}(\mathcal{H}) \leq \epsilon^U(\pi)$ is true with confidence $1 - \beta$





# $\epsilon_L(\pi) \leq \mathsf{R}(\mathcal{H}) \leq \epsilon^U(\pi)$ is true with confidence $1 - \beta$





#### Main result (cont'd)

 $\epsilon_L(\pi) \leq \mathsf{R}(\mathcal{H}) \leq \epsilon^U(\pi)$  is true with confidence  $1 - \beta$ 











 $\min_Q$  $\operatorname{volume}(Q)$ s.t.  $p^{(i)} \in Q$ ,  $i = 1, \dots, N$ 





 $\min_{Q} \quad \text{volume}(Q) \\ \text{s.t.} \quad p^{(i)} \in Q, \\ i = 1, \dots, N$ 

# Q inappropriate for p if p remains outside

Risk = mass outside





 $\min_{Q} \quad \text{volume}(Q) \\ \text{s.t.} \quad p^{(i)} \in Q, \\ \quad i = 1, \dots, N$ 

# Q inappropriate for p if p remains outside

Risk = mass outside

**Compression** = vertexes of the convex hull

 $\pi$  = no. of vertexes



#### N = 500, risk assessment via sample compression theory





#### N = 500, risk assessment via sample compression theory





#### *N* = 500, risk assessment via test dataset (new 500 scenarios)





#### N = 1000, risk assessment via sample compression theory























Many learning algorithms (SV methods, all scenario optimization schemes...) naturally satisfy compression properties... many yet to be discovered...



However, many others do not... notably: Neural Networks





Many learning algorithms (SV methods, all scenario optimization schemes...) naturally satisfy compression properties... many yet to be discovered...



However, many others do not... notably: Neural Networks

**Idea** – the Pick-to-Learn (P2L) algorithm:

construct a meta-algorithm that builds on the existing learning algorithm as a block-box to induce the compression properties



INPUT: scenarios  $\delta_1, \delta_2, \ldots, \delta_N$ , learning algorithm  $\mathfrak{L}$ , initial decision  $\mathcal{H}_0$ 

Initialization: 
$$T = \emptyset, V = (\delta_1, \dots, \delta_N), \mathcal{H} = \mathcal{H}_0$$



 $\overline{\delta}$  = element in V for which  $\mathcal{H}$  is least appropriate



**P2L:** 
$$\delta_1, \ldots, \delta_N \to \mathcal{H}$$

 $\implies$  new learning-based decision scheme  $\mathfrak{L}'$ 

**P2L:** 
$$\delta_1, \ldots, \delta_N \to T$$

 $\Rightarrow$  compression function  $\kappa'$  associated to  $\mathfrak{L}'$ 



**P2L:** 
$$\delta_1, \ldots, \delta_N \to \mathcal{H}$$

 $\implies$  new learning-based decision scheme  $\mathfrak{L}'$ 

**P2L:** 
$$\delta_1, \ldots, \delta_N \to T$$

$$\Rightarrow$$
 compression function  $\kappa'$  associated to  $\mathfrak{L}'$ 

**Theorem** (with D. Paccagnan and M. Campi)

Preference and coherence hold true!

the risk of  $\mathcal{H} = \mathfrak{L}'(\delta_1, \dots, \delta_N)$  can be assessed via the size of T

# Numerical example: classification of MNIST dataset

binary digit classification –  $\mathfrak{L} = SGD$  for NN

P2L vs test-set validation (TSV)





train / initialization portion

# **Numerical example: non-linear regression**



P2L vs test-set validation (TSV)

0.4

0.3

0.2

0.1

0





--- risk bound P2L --- risk bound TSV actual risk P2L actual risk TSV



FUTURE AI RESEARCH

Many thanks to prof. Nicolò Cesa-Bianchi for insightful and inspiring discussions

This research is supported by FAIR (Future Artificial Intelligence Research) project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme (M4C2, Investment 1.3, Line on Artificial Intelligence)



#### Relevant articles

- M.C. Campi, S. Garatti. Compression, Generalization and Learning. Journal of Machine Learning Research, 24(339):1-74, 2023.
- D. Paccagnan, M.C. Campi, S. Garatti, The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance. In: Advances in Neural Information Processing Systems 36 (NeurIPS 2023), 2023.



Many thanks to prof. Nicolò Cesa-Bianchi for insightful and inspiring discussions



Relevant a

- M.C. Campi, S. Garatti. Compression, Generalization and Learning. Journal of Machine Learning Research, 24(339):1-74, 2023.
- D. Paccagnan, M.C. Campi, S. Garatti, The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance. In: Advances in Neural Information Processing Systems 36 (NeurIPS 2023), 2023.





$$\min_{Q} \quad \text{volume}(Q) + \rho \sum_{i=1}^{N} \xi_{i}$$
  
s.t. 
$$\operatorname{dist}(p^{(i)}, Q) \leq \xi_{i}$$
$$i = 1, \dots, N$$





$$\min_{Q} \quad \text{volume}(Q) + \rho \sum_{i=1}^{N} \xi_{i}$$
  
s.t. 
$$\operatorname{dist}(p^{(i)}, Q) \leq \xi_{i}$$
$$i = 1, \dots, N$$

# Q inappropriate for p if p remains outside

Risk = mass outside



s.t.  $\operatorname{dist}(p^{(i)}, Q) \le \xi_i$  $i = 1, \ldots, N$ 

> Q inappropriate for p if p remains outside

Risk = mass outside

**Compression** = vertexes of the convex hull + violated  $\pi$  = no. of vertexes + violated



i=1

