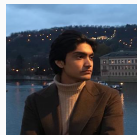




Optimal Continual Learning (CL) has Perfect Memory and is NP-hard

Jeremias Knoblauch^{1,2}, Hisham Husain^{3,4},
Tom Diethe⁵

November 6, 2020



¹University of Warwick, Department of Statistics

²The Alan Turing Institute for Data Science and AI

³Australian National University, College of Engineering & Computer Science

⁴Data61, Sydney

⁵Amazon

Theory:

- (1) Connects set theory to *catastrophic forgetting* in CL
- (2) Avoiding catastrophic forgetting (= optimal CL)
 - (A) is NP-HARD;
 - (B) needs perfect memory.

Practical ramifications:

- (A) CL algorithms = heuristics for NP-HARD problem
- (B) CL with memorization $>$ CL with regularization

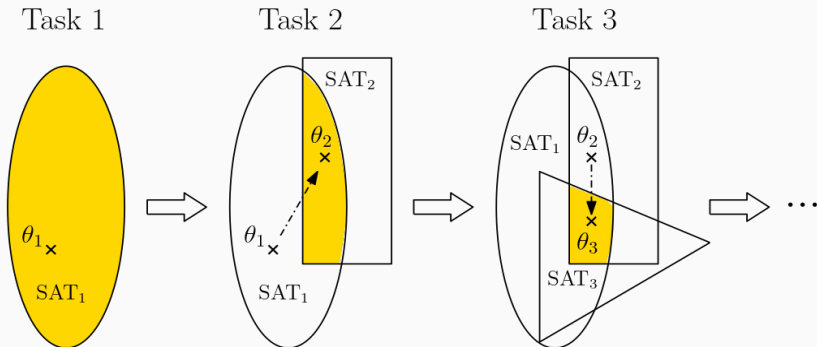
Spotlight: CL and set theory

No catastrophic forgetting = all tasks Satisfy optimality criterion \mathcal{C}

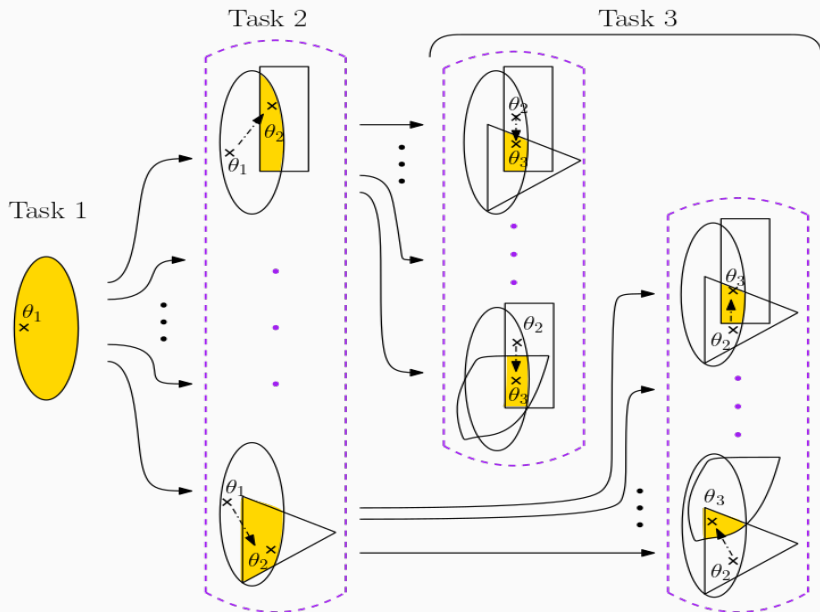
Interpretations using Satisfiability sets:

(1) $\text{SAT}_t = \{\theta \in \Theta : \mathcal{C}(\theta) \text{ is satisfied on task } t\}$.

(2) **No catastrophic forgetting** (= optimal CL) $\iff \theta_t \in \bigcap_{i=1}^t \text{SAT}_i$.

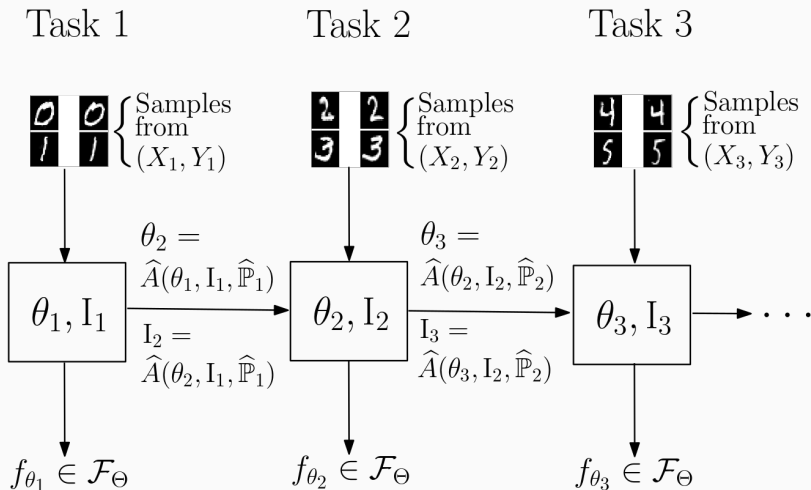


Spotlight: Hardness and memory for CL



- (1) Set Theory & CL
 - (1.1) CL & Catastrophic Forgetting
 - (1.2) Analyzing CL via sets
- (2) Optimal CL: NP-hardness
 - (2.1) The set intersection problem
 - (2.2) NP-hardness results
 - (2.3) A linear model example
- (3) Optimal CL: Perfect Memory
 - (3.1) Defining Perfect Memory
 - (3.2) Memory requirements of CL
 - (3.3) A linear model example
- (4) Practical Ramifications: Memorization vs Regularization

(1.1) CL and Catastrophic Forgetting



Catastrophic Forgetting: θ_3 fails to correctly label Task 1, i.e.



(1.2) Catastrophic Forgetting & Optimality

Notation:

θ_t = parameter value after task t

\mathcal{Q} = empirical distributions of all possible tasks

$\hat{\mathbb{P}}_t \in \mathcal{Q}$ = t -th task's empirical distribution:

$$\hat{\mathbb{P}}_t(x, y) = \frac{1}{n_t} \sum_{i=1}^{n_t} \delta_{(y_i^t, x_i^t)}(y, x), \text{ for } n_t \in \mathbb{N}, y_i^t \in \mathbb{R}, x_i^t \in \mathbb{R}^d$$

\mathcal{C} = Optimality criterion. E.g., linear model & ε -error bound:

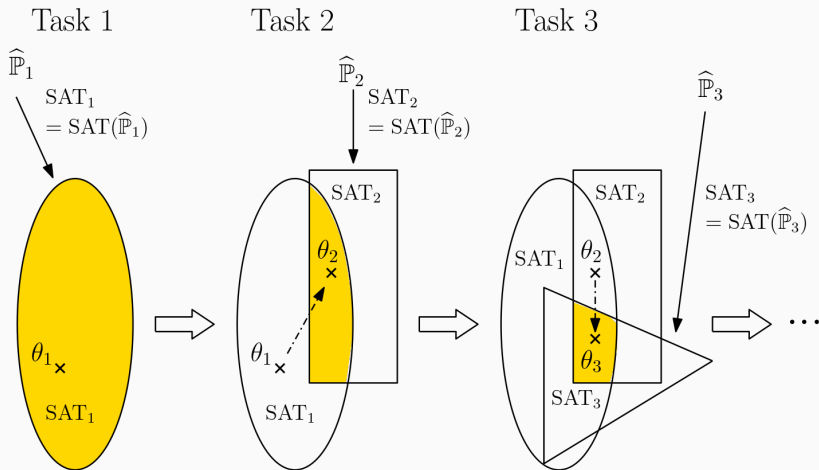
$$\mathcal{C}(\theta, \hat{\mathbb{P}}) = \begin{cases} 1 & \text{if } \frac{1}{n_t} \sum_{i=1}^{n_t} |y_i^t - \theta^T x_i^t| \leq \varepsilon \\ 0 & \text{otherwise.} \end{cases}$$

SAT_t = Satisfiability sets of task t , $\text{SAT}_t = \{\theta \in \Theta : \mathcal{C}(\theta, \hat{\mathbb{P}}_t) = 1\}$

- (i) Optimality on all tasks = No catastrophic forgetting
- (ii) Optimality $\iff \theta_t \in \bigcap_{i=1}^t \text{SAT}_i$
- (iii) Lemma 1: Analysis of optimal CL via SAT_t valid!

(1.2) Catastrophic Forgetting & Optimality

Meaning of Lemma 1: Optimal CL has a set-theoretic interpretation



(2.1) The set intersection problem

Notation:

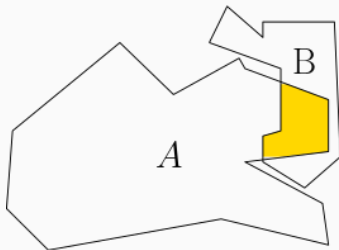
$SAT_Q =$ set of all possible SAT_t

$SAT_{\cap} =$ all finite intersections $\cap_{i=1}^t SAT_i$ of sets in SAT_Q

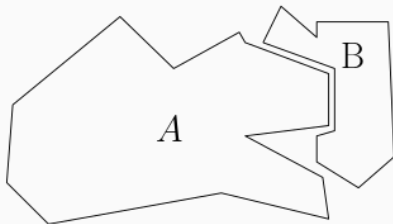
Lemma 2

An optimal CL algorithm is computationally at least as hard as deciding whether $A \cap B = \emptyset$, for $A \in SAT_{\cap}$ and $B \in SAT_Q$.

Situation 1: $A \cap B \neq \emptyset$



Situation 2: $A \cap B = \emptyset$



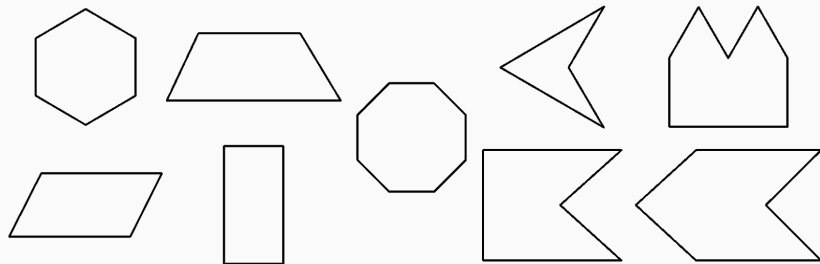
(2.2) NP-Hardness result

Theorem 1

If $\text{SAT}_{\mathcal{Q}} \supseteq S$ or $\text{SAT}_{\cap} \supseteq S$ so that S is the set of tropical hypersurfaces or the set of polytopes on Θ , then optimal CL is NP-HARD.

Proof sketch.

Combine Lemma 2 with the fact that the intersection problem is NP-COMPLETE if S is the set of tropical hypersurfaces or the set of polytopes on Θ . □



(2.3) A linear model example for NP-Hardness

- (i) θ = coefficient vector of linear model
- (ii) \mathcal{Q} = all tasks with empirical measures

$$\hat{\mathbb{P}}_t(x, y) = \frac{1}{n_t} \sum_{i=1}^{n_t} \delta_{(y_i^t, x_i^t)}(y, x), \text{ for } n_t \in \mathbb{N}, y_i^t \in \mathbb{R}, x_i^t \in \mathbb{R}^d$$

- (iii) \mathcal{C} = ε -upper bound criterion in the L_1 -norm:

$$\mathcal{C}(\theta, \hat{\mathbb{P}}) = \begin{cases} 1 & \text{if } \frac{1}{n_t} \sum_{i=1}^{n_t} |y_i^t - \theta^T x_i^t| \leq \varepsilon \\ 0 & \text{otherwise.} \end{cases}$$

$\implies \text{SAT}_{\cap}$ contains *all* polytopes in Θ , since

$$\text{SAT}_t = \{ \theta \in \Theta : \mathbf{y}^t - \theta^T \mathbf{X}^t \leq \varepsilon \cdot n_t \text{ and } \theta^T \mathbf{X}^t - \mathbf{y}^t \geq \varepsilon \cdot n_t \}$$

(2.3) A linear model example for NP-Hardness

Q: So what? Who cares about simple Linear Models?

Corollary 1

If deciding whether $A \cap B = \emptyset$ for $A, B \in S$ is computationally at least as hard as for the collection of polytopes in Θ , optimal CL is NP-HARD.

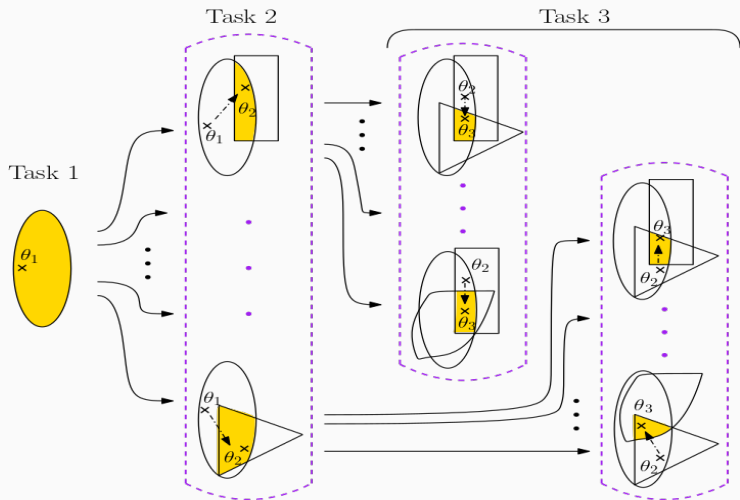
Meaning of Corollary 1:

Complicated nonlinear models / Neural Networks also NP-HARD

(3.1) Defining Perfect Memory

Intuition: We could encounter *any* $\text{SAT}_t \in \text{SAT}_Q$ at task t

\implies Need all information in $\cap_{i=1}^{t-1} \text{SAT}_i$ from tasks $1, \dots, t-1$



(3.1) Defining Perfect Memory

Q: Smallest necessary informational content in $\cap_{i=1}^t \text{SAT}_i$?

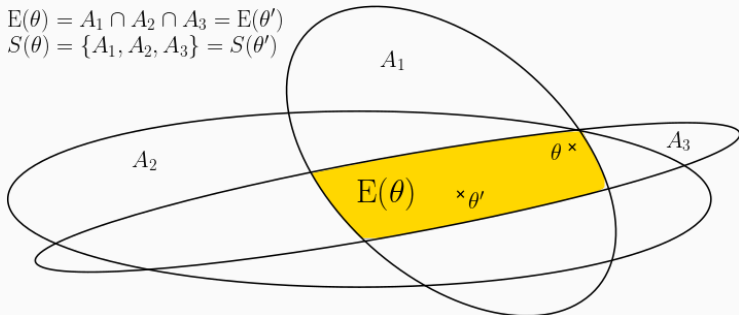
\implies The subset of $\cap_{i=1}^t \text{SAT}_i$ without *equivalent/redundant* elements!

Definition 1 (Equivalence set (\approx 'redundance class'))

For $\theta \in \Theta$, define $S(\theta) = \{A \in \text{SAT}_{\mathcal{Q}} : \theta \in A\}$ and the equivalence sets

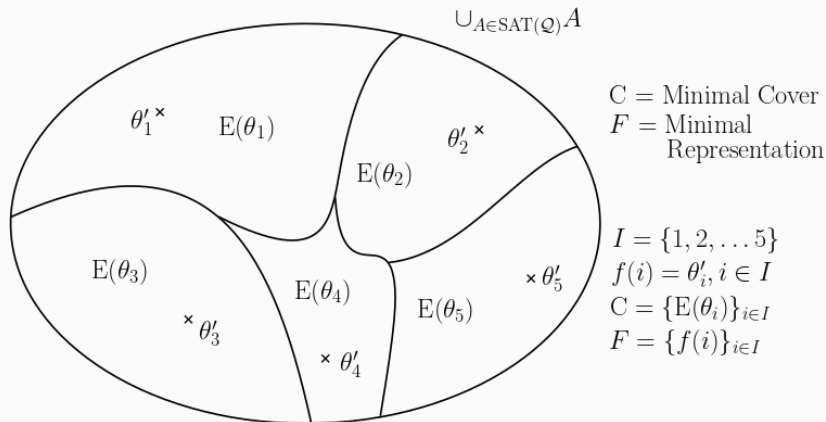
$$E(\theta) = \bigcap_{A \in S(\theta)} A.$$

$$\begin{aligned} E(\theta) &= A_1 \cap A_2 \cap A_3 = E(\theta') \\ S(\theta) &= \{A_1, A_2, A_3\} = S(\theta') \end{aligned}$$



(3.1) Defining Perfect Memory

Perfect Memory = store element of each $E(\theta)$ s.t. $E(\theta) \subseteq \bigcap_{i=1}^t \text{SAT}_i$



(3.2) CL Memory Requirements

Lemma 4

For optimal CL algorithms, there exists $h : \Theta \times i \rightarrow 2^\Theta$ for which $h(\theta_t, I_t) = C_t$ is s.t. $C_t \cap A = \emptyset \iff \bigcap_{i=1}^t \text{SAT}_i \cap A = \emptyset, \forall A \in \text{SAT}_Q$.

Meaning of Lemma 4:

Optimal CL memorizes enough to solve set intersection problem

Corollary 2

If \mathcal{C} and \mathcal{Q} are s.t. $E(\theta) \in \text{SAT}_Q$ for all $\theta \in \Theta$, any optimal CL algorithm has perfect memory.

Meaning of Corollary 2:

You can only solve the set intersection problem with perfect memory

(3.3) A linear model example for Perfect Memory

- (i) θ = coefficient vector of linear model
- (ii) \mathcal{Q} = all tasks with empirical measures

$$\hat{\mathbb{P}}_t(x, y) = \frac{1}{n_t} \sum_{i=1}^{n_t} \delta_{(y_i^t, x_i^t)}(y, x), \text{ for } n_t \in \mathbb{N}, y_i^t \in \mathbb{R}, x_i^t \in \mathbb{R}^d$$

- (iii) \mathcal{C} = ε -upper bound criterion in the L_1 -norm:

$$\mathcal{C}(\theta, \hat{\mathbb{P}}) = \begin{cases} 1 & \text{if } \frac{1}{n_t} \sum_{i=1}^t |y_i^t - \theta^T x_i^t| \leq \varepsilon \\ 0 & \text{otherwise.} \end{cases}$$

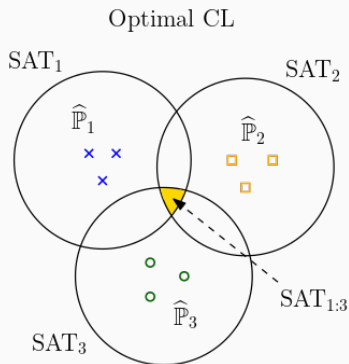
$\implies \text{SAT}_{\mathcal{Q}}$ contains $\{\theta\} = \text{E}(\theta)$ for all $\theta \in \Theta$.

(Take $\mathbf{X}^t, \mathbf{y}^t$ s.t. $\mathbf{y}^t - \varepsilon n_t = \theta^T \mathbf{X}^t$ is solved by **unique** $\theta \in \Theta$)

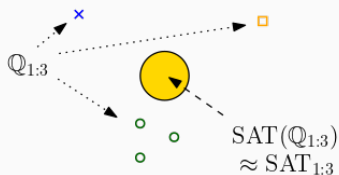
(4) Practical Ramifications: Memorization vs Regularization

Implications for algorithm design:

- (1) Optimal CL = moving $\cap_{i=1}^{t-1} \text{SAT}_i \longrightarrow \cap_{i=1}^{t-1} \text{SAT}_i \cap \text{SAT}_t$
- (2) We need perfect memory to do that without errors
 \Rightarrow memorization should work better than regularization



CL with Core Sets/Replay/Memory



Theory:

- (1) Connects set theory to *catastrophic forgetting* in CL
- (2) Avoiding catastrophic forgetting (= optimal CL)
 - (A) is NP-HARD;
 - (B) needs perfect memory.

Practical ramifications:

- (A) CL algorithms = heuristics for NP-HARD problem
- (B) CL with memorization $>$ CL with regularization

Contact me/follow my research at

E-mail address: j.knoblauch@warwick.ac.uk

Personal Website: <https://jeremiasknoblauch.github.io/>

Twitter Handle: [@LauchLab](https://twitter.com/LauchLab)

