



Introduction

Given: Dependent (β -mixing) data $Z_{1:n} = (X, Y)_{1:n}$ and a hypothesis class \mathcal{F} **Output:** a predictor $\hat{f} \in \mathcal{F}$ that "minimizes":

$$\mathbf{ER}(\hat{f}) \triangleq \mathbf{E}_{X,Y} \| \hat{f}(X) - Y \|^2 - \min_{f \in \mathcal{F}} \mathbf{E}_{X,Y} \| \| f(X) - Y \|^2$$

Main: establish instance optimal rates for ERM

Distinguish between:

Realizability: $Y_i = f_{\star}(X_i) + W_i$, for some MDS "white noise" $W_{1:n}$ **Agnostic:** no general relation between X and $Y \Rightarrow$ life is significantly harder

• Can still define $W_{1:n}$ via $f_{\star} \in \operatorname{argmin}_{f \in \mathscr{F}} E_{X,Y} ||f(X) - Y||^2$ and $W_i \triangleq Y_i - f_{\star}(X_i)$

Motivation



IID learning is very well understood — rich theory:

 Well-known asymptotics, sharp rates for ERM, lots of algorithmic results Key issue: temporal dependence in data not allowed in IID learning!



Goal: can we build a sharp theory for dependent learning? Problem: Blocking typically deflates the sample size Blocking transforms n dep. samples $m = n/t_{mix}$ independent "blocks" $Z_{1:n} \Rightarrow \tilde{Z}_{1:t_{mix}}, \tilde{Z}_{k+1:2t_{mix}}, \tilde{Z}_{2k+1:3t_{mix}}, \dots$ (*m* independent blocks) Can now apply standard results for independent data to the "blocks" Generically employed, deflates rate of converge by a factor t_{mix}... Classical asymptotics tell us this is not instance-optimal!

Sharp Rates in Dependent Learning Theory: Avoiding Sample Size Deflation for the Square Loss Ingvar Ziemann (Penn), Stephen Tu (USC), George J. Pappas (Penn), Nikolai Matni (Penn)

$$\hat{f} \in \operatorname{argmin}_{f \in \mathscr{F}} \frac{1}{n} \sum_{i=1}^{n} ||f(X_i) - Y_i||^2$$

$$MDS \text{ "white noise" } W_{1:n}$$

Theorem [ZTPM2024]: Let \mathcal{F} be either a ERM for β -mixing stationary data $(X, Y)_{1 \cdot n}$

- $||f||_{\Psi_n} \leq L||f||_{I^2}^{\eta}$ for some $\eta \in (0,1]$ v
- *n* is greater than a polynomial *burn-in* n = n = n = n

We have that with probability at least 1 - d

 $ER(\hat{f}) \leq \sigma^2 \times \frac{CO1}{CO1}$

 $\sigma^{2} \triangleq \lim_{n \to \infty} \sup_{g \in (\mathcal{F} - \mathcal{F}) \cap o(1)S_{L^{2}}} \operatorname{Var}\left(\frac{1}{\sqrt{n}}\sum_{i=1}^{n} \left\langle W_{i} \right\rangle \right)$

We also show that the weak subG class condition holds

- For finite hypothesis classes
- For smooth (in the input) hypothesis classical
- For parametric hypothesis classes (Loja

Dependency Deflation?

Stable GLM, numerical ERM experiment

GLM: $Y_t \triangleq X_{t+1} = \phi(A_{\star}X_t) + W_t$, w/ ϕ a known link function (here leaky relu), $W_t \sim N(0, t)$ Dependency ($\rho = \rho(A_{\star})$) doesn't seem to hurt (at least not for large T = n):



Key takeaway: Under realizability (i.e., $\mathbf{E}[Y_t \mid X_t] = f_{\star}(X_t)$), the leading variance proxy term tends to the same variance as in the independent case In other words: dependence only hurts without realizability and the degradation is graceful

Contribution

a 1) convex	<i>x or 2) realizab</i> se further that:	le class	and let \hat{f}	be the			
μ $\mathbf{v} = \mathbf{v} = \frac{1}{p} \left[\frac{\Delta}{v} - \frac{1}{p} \right] \left[f(\mathbf{v}) \right]$							Firs
wnere II/ II	$\Psi_p = \sup m$ $m \ge 1$	$ = P \prod_{X \in \mathcal{X}} (X) $	II <i>Lm</i>				Defi
<i>in</i> in problem constants (<i>t_{mix}</i> , dim, etc.)							Q_n
δ :							
mplexit	$y(\mathcal{F}) + \log$	$g(1/\delta)$					On t
	n						
$\left(\frac{g(X_i)}{\ g\ _{L^2}} \right)$ rS_{L^2} : rad r sphere in L^2							Defi
Var(W) if strictly realizable)							M_{r}
		$ = + e^{i} = - e^{i} + e^{i} = - e^{i} = - e^{i} + e^{i} = - e^{$					
	I ne term coi	eterm complexity(F):					
	 Is formally 	; formally a local Talagrand γ_2 -functional					
	 Scales like 	$d_{\mathbf{x}}$ for a	$l_{\mathbf{x}}$ -dimens	sional	linear		
asses regression							
ija)	 More gene parametric entropy) 	rally sca classes	les like <i>d</i> (controlle	for <i>d</i> - ad by	dimensi	ional	Insi the The
$(0,\sigma_W^2)$ iid	Overvie	w of Res	sults: dep). data	& squa	re loss	•
	Paper	Hypothesis Class	Rate (Hiding logs)	Mixing/ Stability	Guarantee	Realizability	
	[SMTJR2018]	Linear	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Marginal	Param	Required	
	[SO2020]	Generalized Lin	$t_{\rm mix} n^{-1} \sigma^2 \times {\rm complexity}(\mathcal{F})$	Strict	Param recovery	Required	
	[KNJN2021]	Generalized Lin	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Marginal	Param recovery	Required	
	[RBE2021]	General	$t_{\rm mix} n^{-1+\epsilon} \sigma^2 \times {\rm complexity}(\mathcal{F})$	Strict	Excess Risk	Not Required	l -
've given us:	[ZT 2022]	General (4-2 hypcon)	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Strict	Excess Risk	Required	
omplexity(F)	[SOO2022]	Bilinear	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Marginal	Param Recovery	Required	
<i>n</i>	[ZTPM 2023]	Linear	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Strict	Excess Risk	Not Required	
ove this term	[ZTPM 2024]	General (weak subG)	$n^{-1}\sigma^2 \times \text{complexity}(\mathcal{F})$	Strict	Excess Risk	Not Required	
	Sharp Rat	es* withou	ut Realizabil	ity.	(* up to a	logarithm)	



Proof Strategy

t since \mathscr{F} is convex or realizable $\mathbf{ER}(\hat{f}) \leq \mathbf{E}_{\mathbf{X}} \| \hat{f}(\mathbf{X}) - f_{\star}(\mathbf{X}) \|^2$ ine the quadratic process:

$$(f) \triangleq \mathbf{E}_{\mathbf{X}} \| \| f(\mathbf{X}) - f_{\star}(\mathbf{X}) \|^{2} - \frac{1 + \epsilon}{n} \sum_{i=1}^{n} \| \| f(\mathbf{X}_{i}) - f_{\star}(\mathbf{X}_{i}) \|^{2}$$

the event $\{Q_{n,\epsilon}(f) \leq 0 : \forall f \in \mathscr{F} \setminus rS_{L^2}\}$ we have:

$$\|(X) - f_{\star}(X)\|^{2} \le r + \frac{1 + \epsilon}{rn} \sum_{i=1}^{n} 2(1 - \mathbf{E}') \left\langle W_{i}, \frac{r[\hat{f}(X_{i}) - f_{\star}(X_{i})]}{\|\hat{f}(X_{i}) - f_{\star}(X_{i})\|_{L^{2}}} \right\rangle$$

ine the multiplier process:

$$f_{E}(f) \triangleq \sum_{i=1}^{n} 2(1 - \mathbf{E}') \left\langle W_{i}, f(X_{i}) - f_{\star}(X_{i}) \right\rangle$$

nce: the proof boils down to unif. controlling $Q_{n,\epsilon}(f)$ and $M_{n,\epsilon}(f)$ over localized class $(\mathcal{F} - \mathcal{F}) \cap rS_{L^2}$ and:

complexity(\mathcal{F}): soln. to $r \asymp$ sup $M_{n,\epsilon}(f)$ $f \in (\mathcal{F} - \mathcal{F}) \cap rS_{I2}$

ght from Mendelson (2014), $Q_{n,\epsilon}(f)$ does not affect the rate, but only burn-in \Rightarrow can freely block to control $Q_{n_c}(f)$

- e crux is to control $M_{n,c}(f)$ carefully. We combine:
- A version of Bernstein's inequality (Maurer and Pontil 2021)

For
$$V_i(f) \triangleq (1 - \mathbf{E}') \left\langle W_i, f(X_i) - f_{\star}(X_i) \right\rangle$$
 with $||f||_{L^2} = r$:

$$\frac{1}{n} V_i \lesssim 2\sqrt{\frac{\mathbf{V}(\bar{V})\ln(1/\delta)}{n}} + \frac{t_{\min} \|V\|_{\Psi_p} \log(1/\delta)}{n} \qquad \bar{V} = \frac{1}{t_{\min}} \sum_{i=1}^k V_i$$

Localization ($\|I\|_{L^2} = r$)

- Balances both RHS terms \Rightarrow reintroduces the mixing time in the rate The weak subG class condition
- Breaks the balance and makes the variance term dominant again • Tail bounds via generic chaining (Dirksen 2015) for local unif. control