

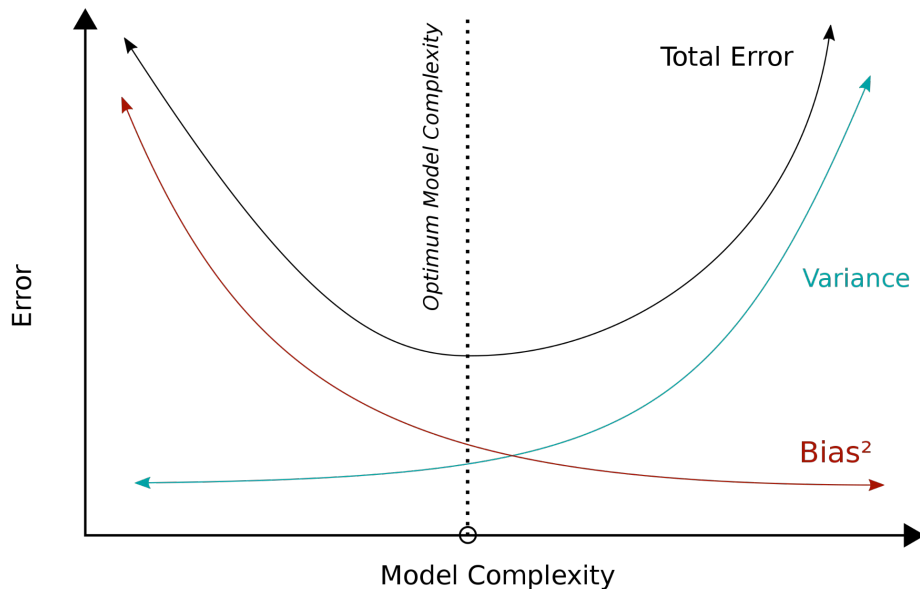
Double Descent



Данил
Шешеня
213

Bias-variance trade-off

- Теоретический концепт
- С ростом сложности модели падает смещение, но растет дисперсия



Житейская мудрость

- Не делать слишком сложные модели
- Тем не менее, большие модели существуют и хорошо работают
- Делать сложные модели!!
- Оба подхода верны

Effective Model Complexity (EMC)

Максимальное число элементов в обучающей выборке, при котором модель достигает ~ 0 ошибки

Пусть:

- $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ – выборка из распределения \mathcal{D}
- $\mathcal{T}(S)$ – процедура обучения
- $\varepsilon > 0$ (на практике ≈ 0.1)
- $\text{Error}_S(M)$ – средняя ошибка модели M на обучающей выборке S

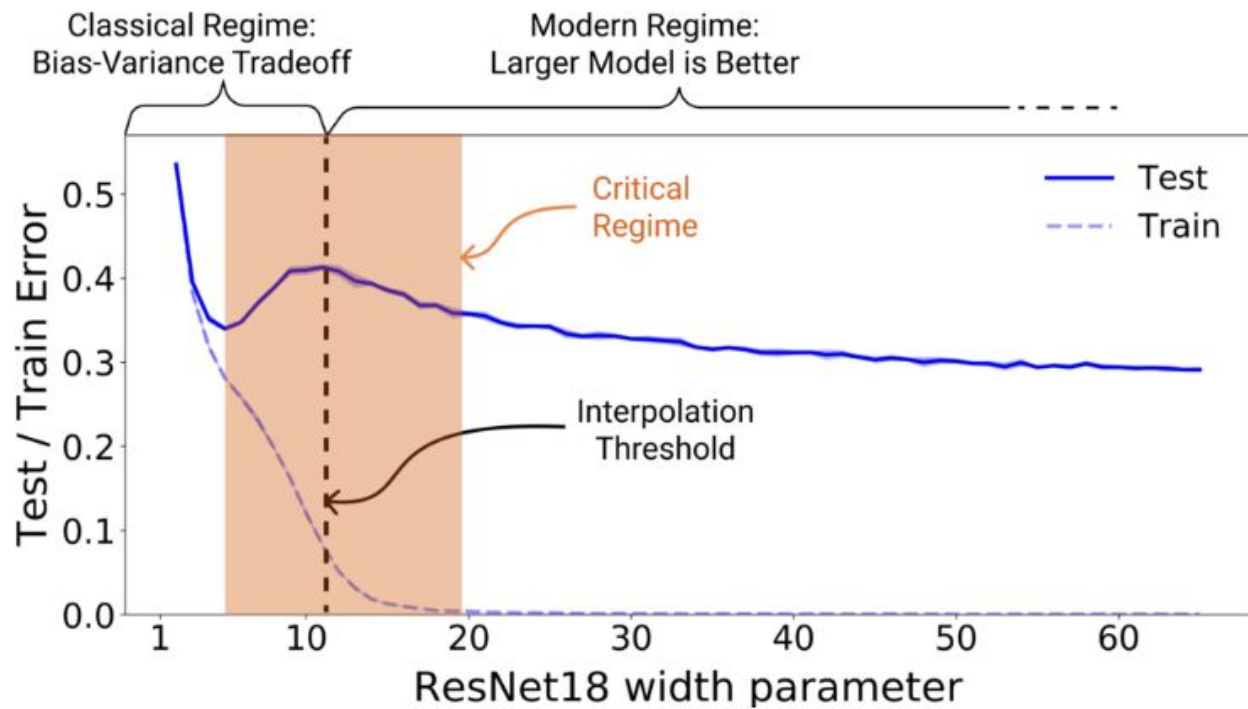
$$\text{EMC}_{\mathcal{D}, \varepsilon}(\mathcal{T}) := \max\{n \mid \mathbb{E}_{S \sim \mathcal{D}^n}[\text{Error}_S(\mathcal{T}(S))] \leq \varepsilon\}$$

Зависит от распределения!

Гипотеза

Для любых \mathcal{D} , \mathcal{T} , $\varepsilon > 0$, если рассмотреть задачу классификации n элементов выборки

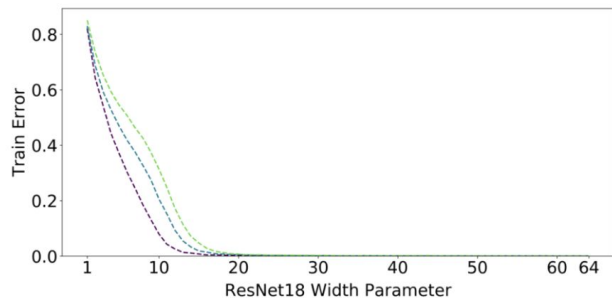
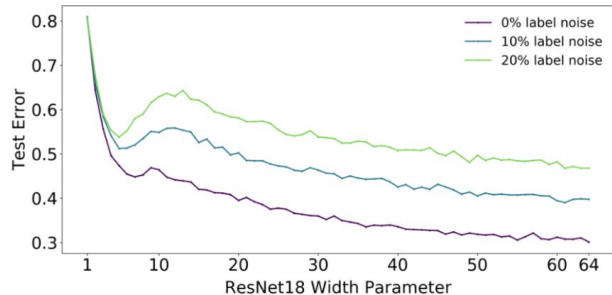
- *Недопараметризованный режим.* Если $EMC \ll n$, то с ростом EMC падает ошибка на тесте
- *Перепараметризованный режим.* Если $EMC \gg n$, то с ростом EMC падает ошибка на тесте
- *Критически параметризованный режим.* Если $EMC \sim n$, то с ростом EMC ошибка на тесте падает или **растет**



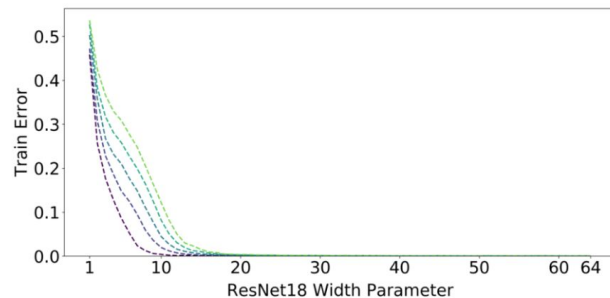
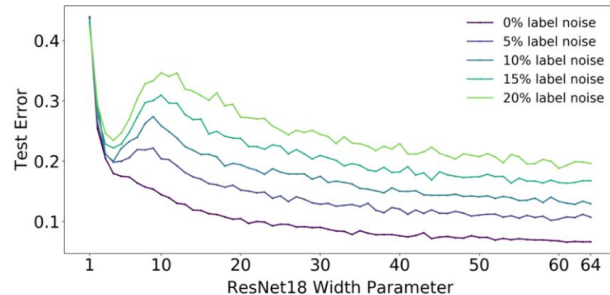
От чего зависит ЕМС?

- Архитектура модели (CNN, ResNet, Transformer)
- Параметры сети (ширина, число слоев и т.д)
- Регуляризация / аугментации / ...
- Число эпох
- Оптимизатор (SGD, Adam)
- Уровень шума в данных

Model-Wise Double Descent

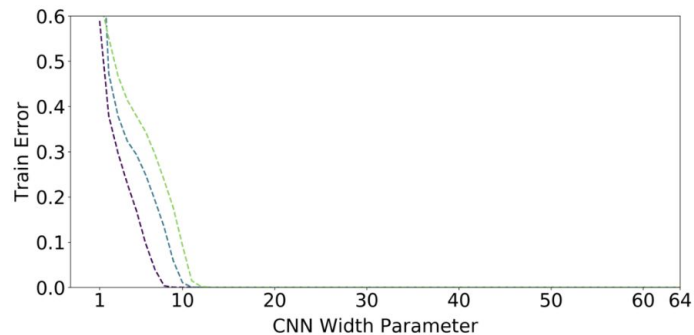
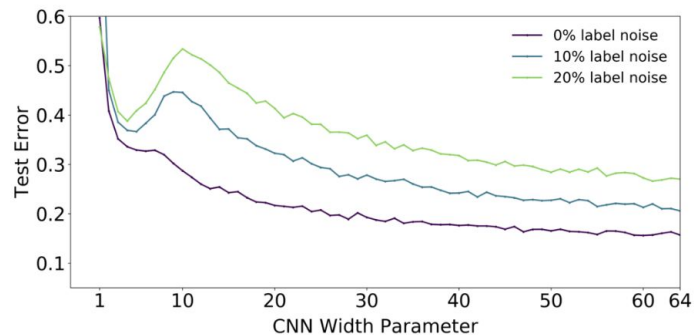


(a) **CIFAR-100.** There is a peak in test error even with no label noise.

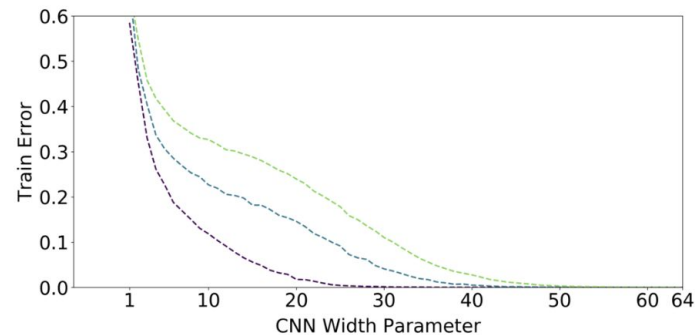
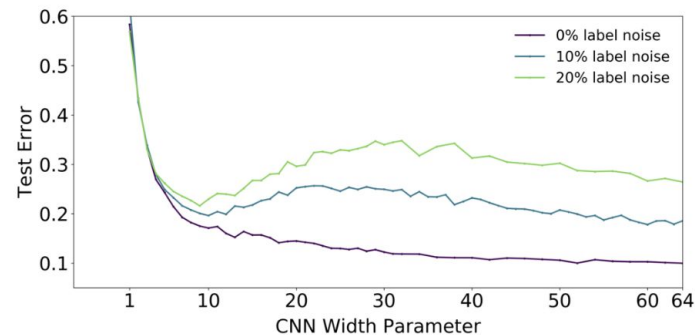


(b) **CIFAR-10.** There is a “plateau” in test error around the interpolation point with no label noise, which develops into a peak for added label noise.

Model-Wise Double Descent

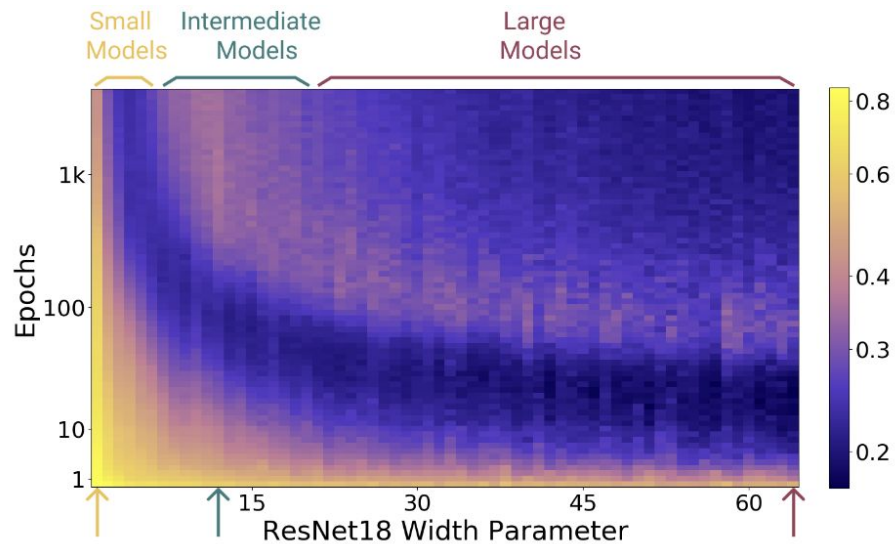
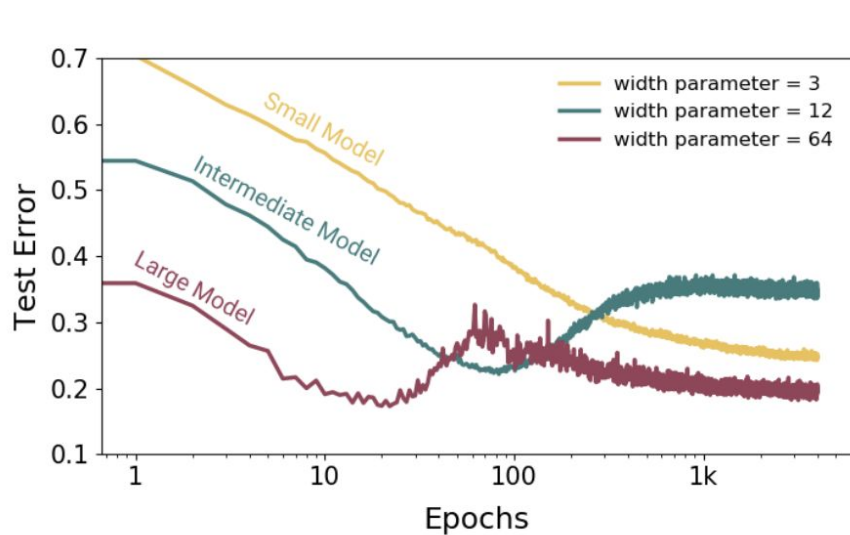


(a) Without data augmentation.



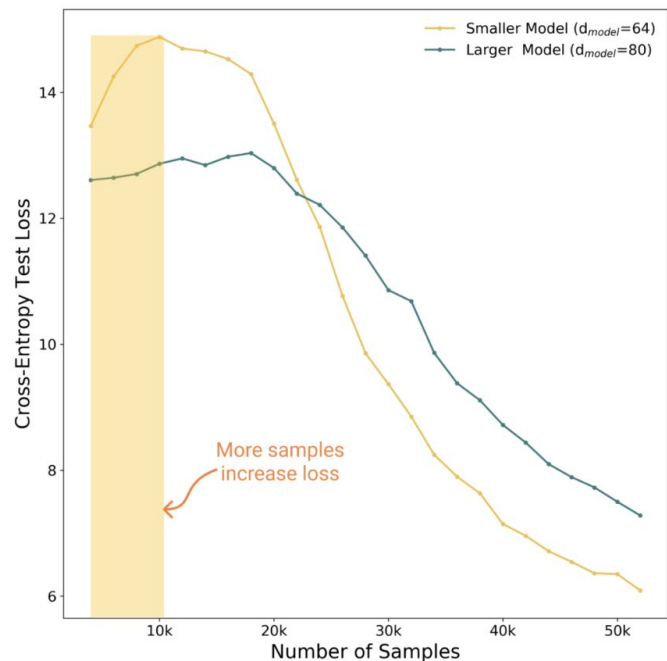
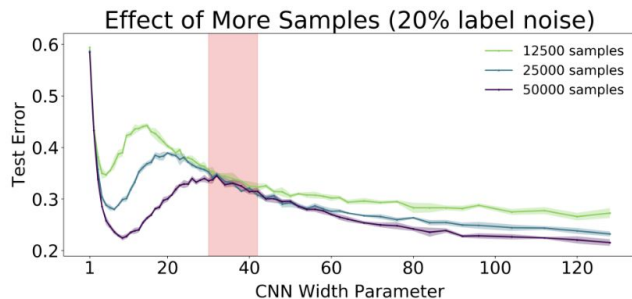
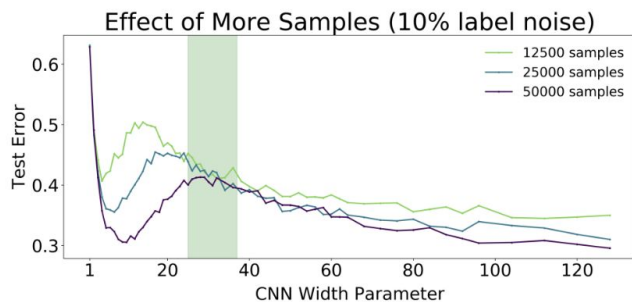
(b) With data augmentation.

Epoch-Wise Double Descent



Sample-Wise Non-Monotonicity

Ситуации $EMC \sim n$ можно достичь меняя n



Результаты

Dataset	Architecture	Opt.	Aug.	% Noise	Double-Descent		Figure(s)		
					Model	Epoch			
CIFAR 10	CNN	SGD	✓	0	✗	✗	5, 27		
			✓	10	✓	✓	5, 27, 6		
			✓	20	✓	✓	5, 27		
				0	✗	✗	5, 25		
				10	✓	✓	5		
				20	✓	✓	5		
		SGD + w.d.	✓	20	✓	✓	21		
			Adam		0	✓	—	25	
			ResNet	Adam	✓	0	✗	✗	4, 10
					✓	5	✓	—	4
	✓	10			✓	✓	4, 10		
	Various	✓		15	✓	✓	4, 2		
		✓		20	✓	✓	4, 9, 10		
		✓		20	—	✓	16, 17, 18		
	(subsampled)	CNN	SGD	✓	10	✓	—	11a	
	(adversarial)	ResNet	SGD	✓	20	✓	—	11a, 12	
				0	Robust err.	—	26		
CIFAR 100	ResNet	Adam	✓	0	✓	✗	4, 19, 10		
			✓	10	✓	✓	4, 10		
			✓	20	✓	✓	4, 10		
	CNN	SGD		0	✓	✗	20		
IWSLT '14 de-en	Transformer	Adam		0	✓	✗	8, 24		
(subsampled)	Transformer	Adam		0	✓	✗	11b, 23		
WMT '14 en-fr	Transformer	Adam		0	✓	✗	8, 24		

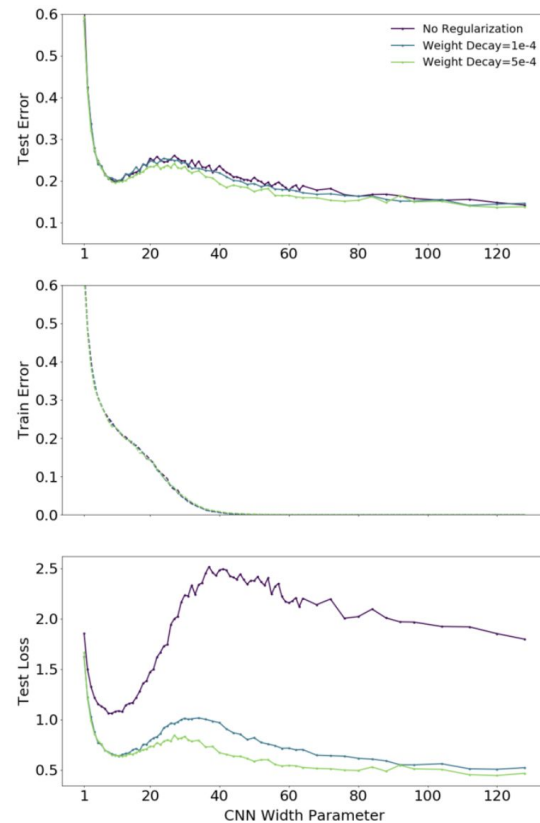
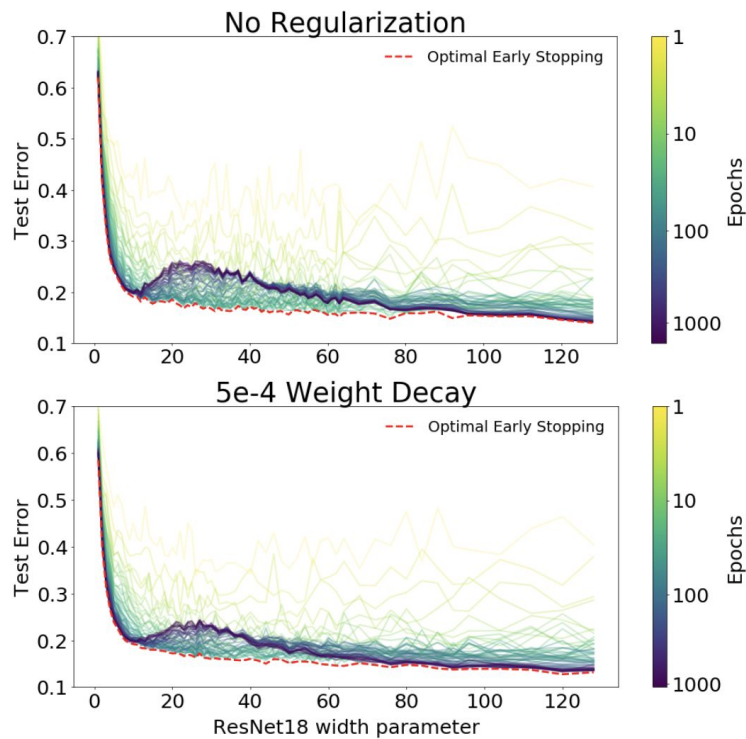
Попытка объяснения эффекта

- По мере увеличения сложности модели, мы приближаемся к интерполяции обучающей выборки
- Такая интерполяционная модель единственна
- Она “шаткая”, и небольшой шум может ее сильно нарушить
- Но дальше есть много моделей, которые интерполируют обучающую выборку и сглаживают шумы

Заключение

- Эффект проявляется довольно устойчиво
- Точные причины эффекта никто не понимает
- Размер критической области тоже никто оценивать не умеет
- Большую роль играет зашумленность данных (model mis-specification)
- Если модель почти запомнила обучающую выборку, то небольшие изменения параметров могут привести к непредвиденным результатам
- Следить за ошибкой на обучающей выборке

Если успею



Еще картинки

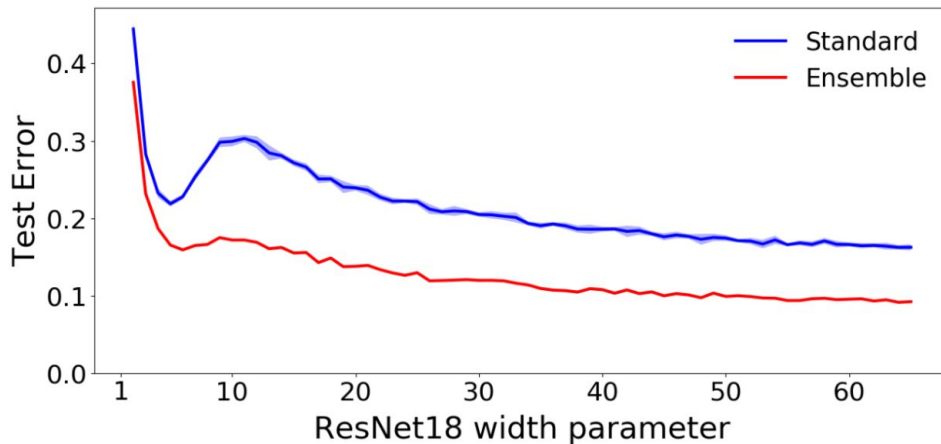


Figure 28: Effect of Ensembling (ResNets, 15% label noise). Test error of an ensemble of 5 models, compared to the base models. The ensembled classifier is determined by plurality vote over the 5 base models. Note that ensembling helps most around the critical regime. All models are ResNet18s trained on CIFAR-10 with 15% label noise, using Adam for 4K epochs (same setting as Figure 1). Test error is measured against the original (not noisy) test set, and each model in the ensemble is trained using a train set with independently-sampled 15% label noise.

Сами разбирайтесь

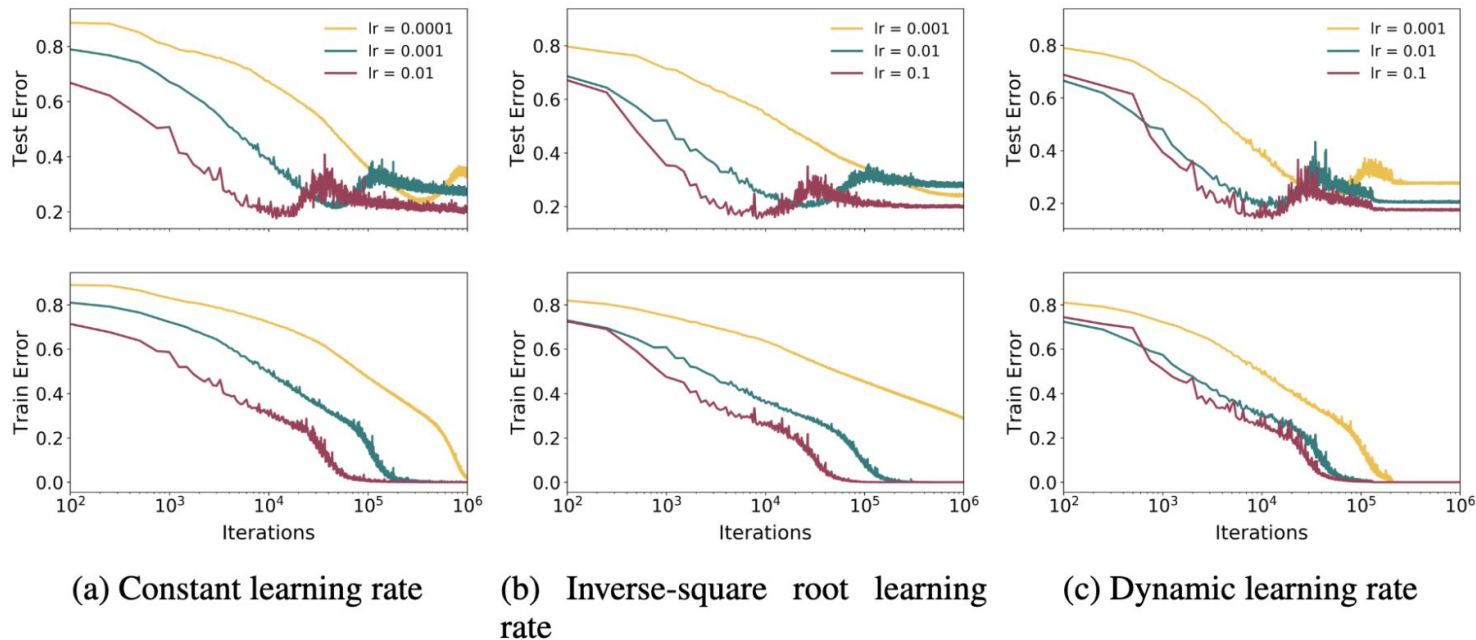
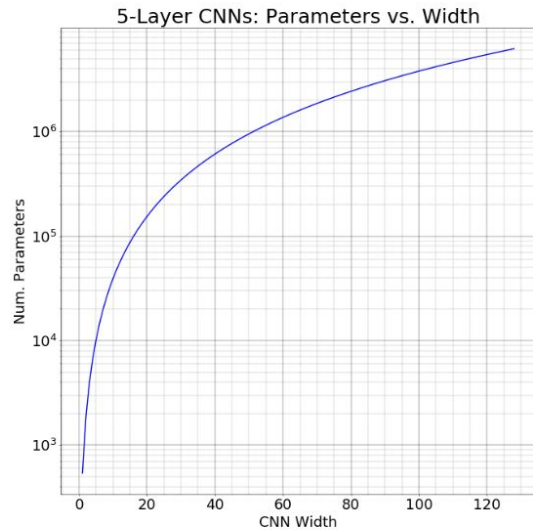
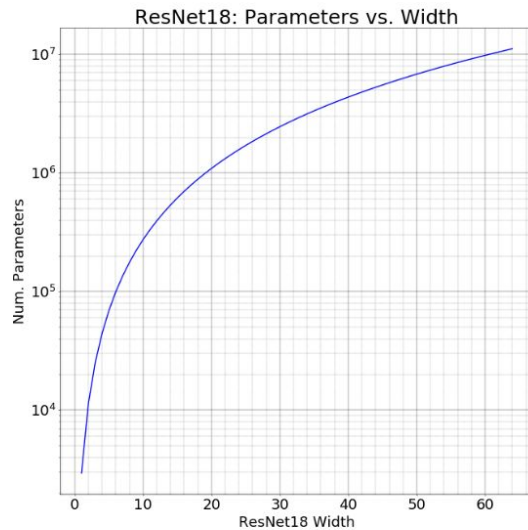


Figure 17: **Epoch-wise double descent** for ResNet18 trained with SGD and multiple learning rate schedules

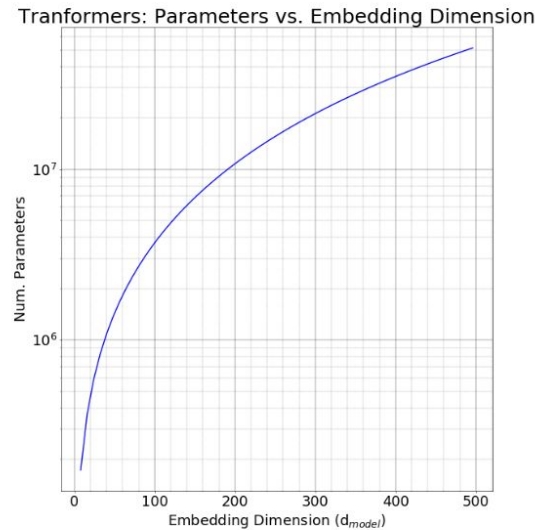
Ты думал я закончил?



(a) 5-layer CNNs



(b) ResNet18s



(c) Transformers