

In-context learning of linear regression

Jingfeng Wu, 10/2023
w/ Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, Peter Bartlett

In-context learning

JW

-0.6615554 => -1.12766266,
0.58568541 => 0.98464498,
-0.06671855 => -0.1269867,
-0.47585609 => -0.80137645,
0.74095414 => ?

```
[>>> x = np.random.randn(5)
[>>> y = x*1.7+np.random.randn(5)*0.01
[>>> x
array([-0.6615554 ,  0.58568541, -0.06671855, -0.47585609,  0.74095414])
[>>> y
array([-1.12766266,  0.98464498, -0.1269867 , -0.80137645,  1.24708834])
>>>
```



It looks like you're trying to find a pattern or rule in a sequence of numbers. Let's analyze the given transformations to see if we can apply the same rule to the number 0.74095414:

Original Number (X) => Transformed Number (Y)

-0.6615554 => -1.12766266
0.58568541 => 0.98464498
-0.06671855 => -0.1269867
-0.47585609 => -0.80137645

It appears that the transformation applied is to multiply the original number by approximately 1.70507279 (rounded to 8 decimal places).

So, to find the transformed number for 0.74095414, you can apply the same rule:

$0.74095414 * 1.70507279 \approx 1.26342623$

So, the transformed number for 0.74095414 is approximately 1.26342623.

Given a pretrained model,
solving new task by

- fine-tuning (backward passes)
- or just a forward pass (ICL)

Linear regression with a Gaussian prior

$$X \in \mathbb{R}^{* \times d}, \quad Y \in \mathbb{R}^*, \quad x \in \mathbb{R}^d, \quad y \in \mathbb{R}$$

1. draw **task** parameter: $\beta \sim \mathcal{N}(0, \psi^2 I_d)$
2. draw covariate-response: $x \sim \mathcal{N}(0, H), y \sim \mathcal{N}(\beta^\top x, \sigma^2)$
3. draw context examples: each row of $(X, Y) \sim \text{iid} \sim (x^\top, y)$

(ψ^2, σ^2, H) are fixed (determining the meta-task)

context length can vary

Model for (weak) ICL

$$f: \mathbb{R}^{* \times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

example 2: single layer attention

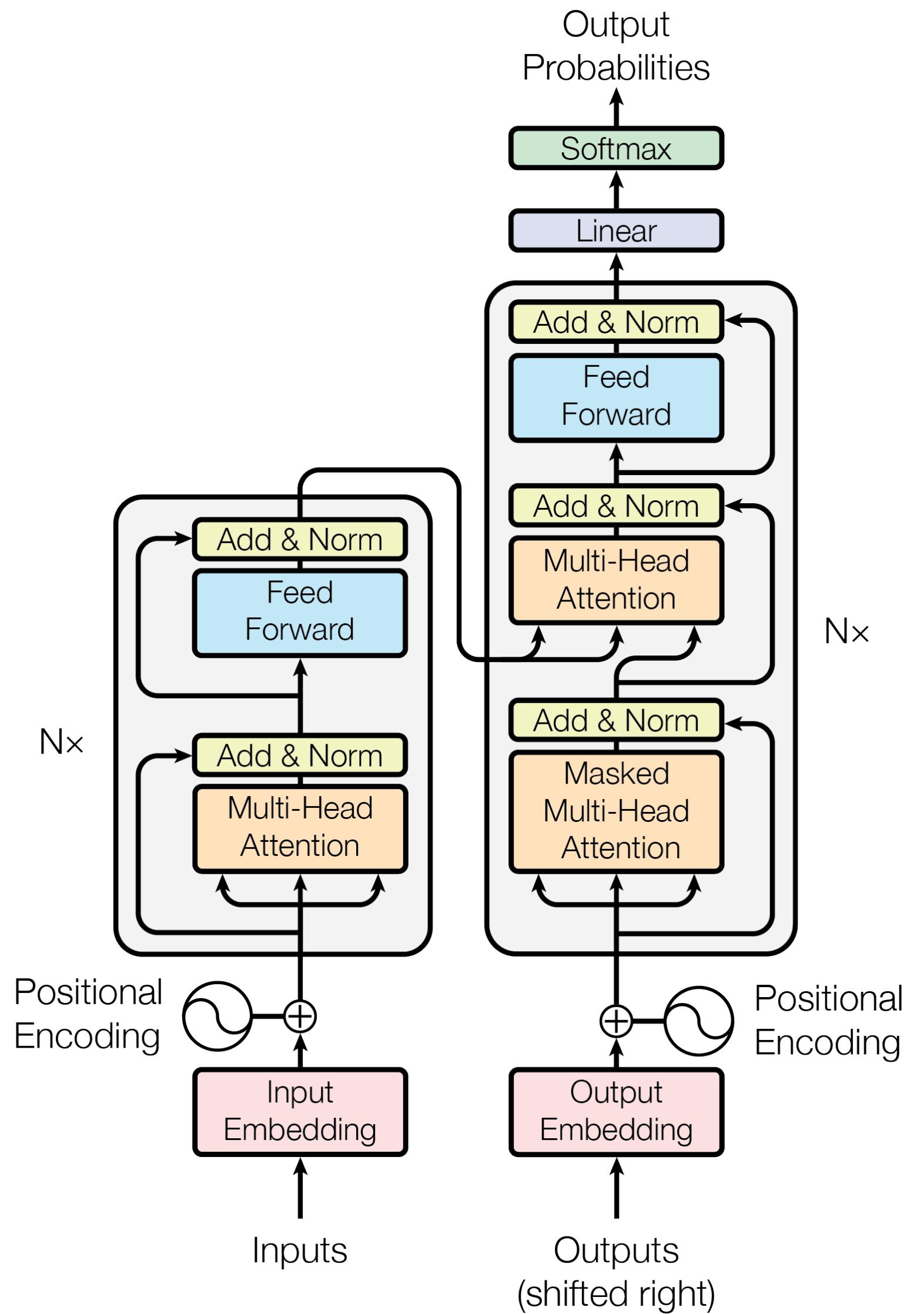
$$Z = \begin{pmatrix} X^\top & x \\ Y^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (n+1)},$$

$$\hat{y} = \left(Z + VZ \cdot \text{sfmx}\left((QZ)^\top(KZ)\right) \right)_{d+1, n+1}$$

V, Q, K are trainable matrix params

example 3: ERM (no param to pretrain)

example 1: transformer



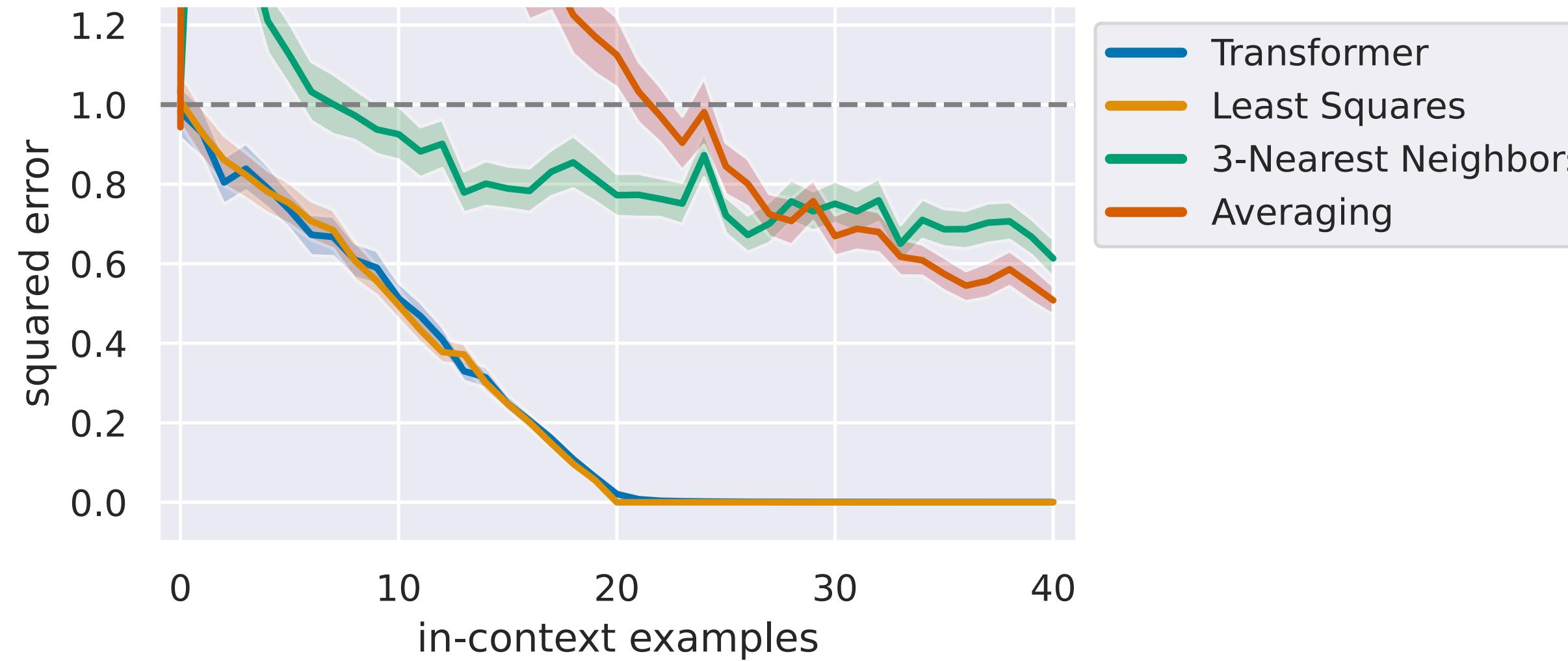
Evidence of ICL by a pretrained transformer

Fixing model f , its **ICL risk at context length n** is

$$f: \mathbb{R}^{* \times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

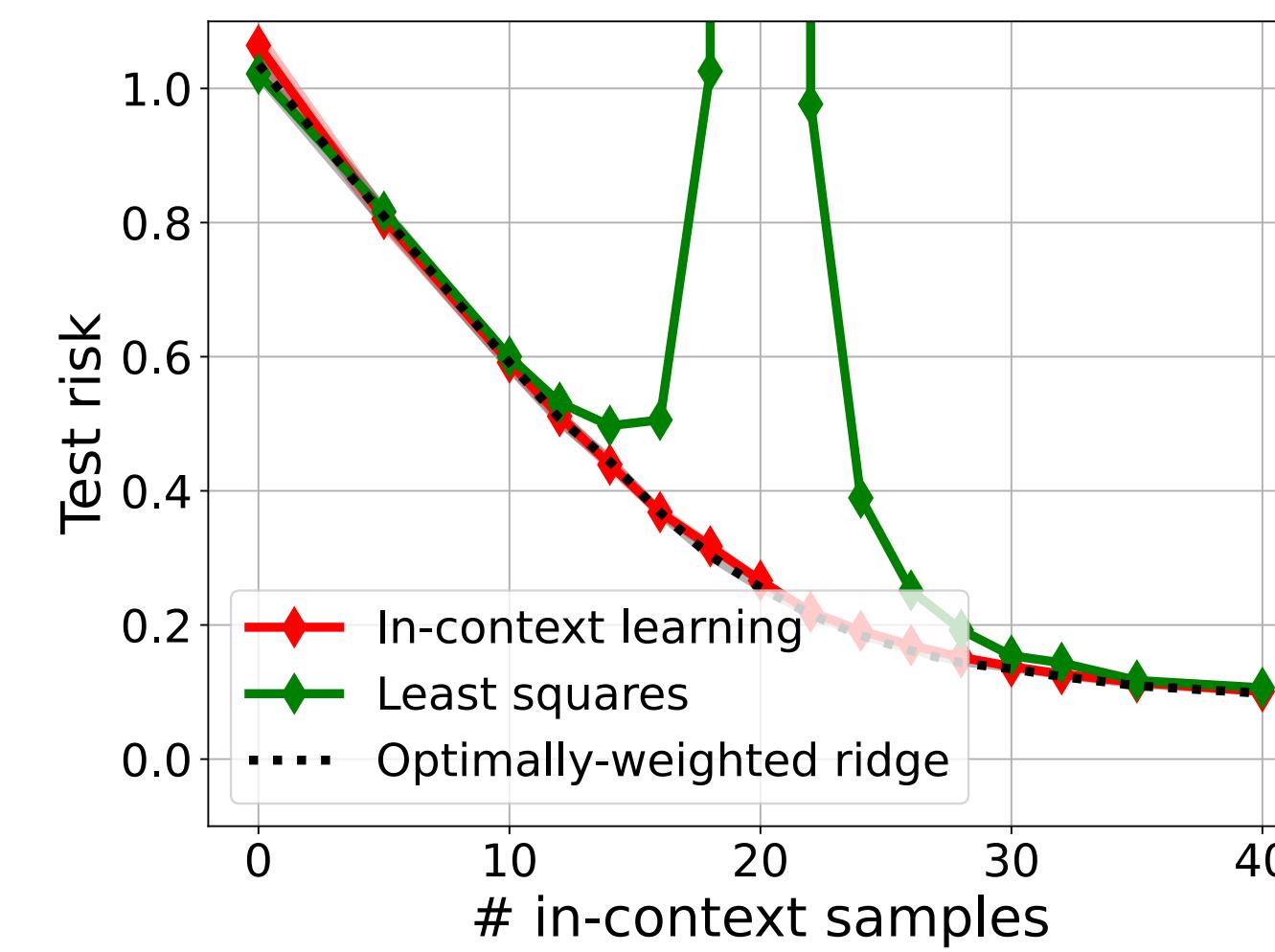
$$(X, Y, x) \mapsto \hat{y}$$

$$R_n(f) = \mathbb{E}_{X, Y, x, v} (f(X, Y, x) - y)^2, \text{ where } \dim(Y) = n$$



$$H = I_{20}, \sigma^2 = 0, \psi^2 = 1$$

$$\hat{y}_{\text{ols}} = \langle (X^\top X)^{-1} X^\top Y, x \rangle$$



$$H = I_{20}, \sigma^2 = 1, \psi^2 = 1$$

$$\hat{y}_{\text{ridge}} = \langle (X^\top X + \sigma^2/\psi^2 \cdot I)^{-1} X^\top Y, x \rangle$$

Garg, Shivam, Dimitris Tsipras, Percy S. Liang, and Gregory Valiant. "What can transformers learn in-context? a case study of simple function classes." *Advances in Neural Information Processing Systems* 35 (2022): 30583-30598.

Akyürek, Ekin, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. "What learning algorithm is in-context learning? investigations with linear models." *arXiv preprint arXiv:2211.15661* (2022).

Li, Yingcong, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. "Transformers as algorithms: Generalization and stability in in-context learning." (2023).

Online pretraining

$$f: \mathbb{R}^{* \times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$
$$(X, Y, x) \mapsto \hat{y}$$

repeat:

1. draw context length $n \in \{1, \dots, N\}$
2. draw new dataset of size $n + 1$: $X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^n, x \in \mathbb{R}^d, y \in \mathbb{R}$
3. update: $f \leftarrow f - \gamma \cdot \nabla_f (f(X, Y, x) - y)^2$

in their experiments: autoregressive loss, ADAM instead of SGD

Garg, Shivam, Dimitris Tsipras, Percy S. Liang, and Gregory Valiant. "What can transformers learn in-context? a case study of simple function classes." Advances in Neural Information Processing Systems 35 (2022): 30583-30598.

Akyürek, Ekin, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. "What learning algorithm is in-context learning? investigations with linear models." *arXiv preprint arXiv:2211.15661* (2022).

Li, Yingcong, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. "Transformers as algorithms: Generalization and stability in in-context learning." (2023).

Some theory about ICL

$$f: \mathbb{R}^{*\times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

Simplification 1: linear attention

single layer attention

$$\hat{y} = \left(Z + VZ \cdot \text{sfmx}\left((QZ)^\top(KZ)\right) \right)_{d+1,n+1}$$

$$Z = \begin{pmatrix} X^\top & x \\ Y^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (n+1)}$$

single layer **linear** attention

$$\hat{y} = \left(Z + VZ \cdot \frac{(QZ)^\top(KZ)}{n} \right)_{d+1,n+1}$$

Ahn, Kwangjun, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. "Transformers learn to implement preconditioned gradient descent for in-context learning." *arXiv preprint arXiv:2306.00297* (2023).

Zhang, Ruiqi, Spencer Frei, and Peter L. Bartlett. "Trained Transformers Learn Linear Models In-Context." *arXiv preprint arXiv:2306.09927* (2023).

$$f: \mathbb{R}^{*\times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

Simplification 2: reparameterization

single layer linear attention

$$\hat{y} = \left(Z + VZ \cdot \frac{(QZ)^T(KZ)}{n} \right)_{d+1,n+1}$$

$$Z = \begin{pmatrix} X^T & x \\ Y^T & 0 \end{pmatrix} \in \mathbb{R}^{(d+1)\times(n+1)}$$

equals to $\hat{y} = \left\langle \boxed{(vW^T)} \cdot \frac{X^T Y}{n}, x \right\rangle$

replace by $\Gamma \in \mathbb{R}^{d \times d}$

if the bottom left $1 \times d$ blocks in V and QK^T are fixed to zeros:

$$V = \begin{pmatrix} * & * \\ 0 & v \end{pmatrix} \text{ and } QK^T = \begin{pmatrix} W & * \\ 0 & * \end{pmatrix}$$

Ahn, Kwangjun, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. "Transformers learn to implement preconditioned gradient descent for in-context learning." *arXiv preprint arXiv:2306.00297* (2023).

Zhang, Ruiqi, Spencer Frei, and Peter L. Bartlett. "Trained Transformers Learn Linear Models In-Context." *arXiv preprint arXiv:2306.09927* (2023).

$$f: \mathbb{R}^{*\times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

\approx one-step GD

single layer linear attention with re-parameterization

$$\hat{y} = \left\langle \Gamma \cdot \frac{X^\top Y}{n}, x \right\rangle, \Gamma \in \mathbb{R}^{d \times d}$$

one-step GD with $w_0 = 0$ and trainable matrix stepsize

$$\hat{y} = \langle \hat{w}, x \rangle, \hat{w} = w_0 - \Gamma \cdot \frac{1}{n} X^\top (Xw_0 - Y)$$

intuition: attention \approx kernel \approx Gram matrix

Ahn, Kwangjun, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. "Transformers learn to implement preconditioned gradient descent for in-context learning." *arXiv preprint arXiv:2306.00297* (2023).

Zhang, Ruiqi, Spencer Frei, and Peter L. Bartlett. "Trained Transformers Learn Linear Models In-Context." *arXiv preprint arXiv:2306.09927* (2023).

Simplification 3: fixed context length

repeat:

1. fix context length $n = N$
2. draw new dataset of size $N + 1$: $X \in \mathbb{R}^{N \times d}, Y \in \mathbb{R}^N, x \in \mathbb{R}^d, y \in \mathbb{R}$
3. update: $\Gamma \leftarrow \Gamma - \gamma \cdot \nabla_{\Gamma} (\hat{y} - y)^2$
$$\hat{y} = \left\langle \Gamma \cdot \frac{X^T Y}{N}, x \right\rangle$$

pretraining risk (i.e., ICL risk **at context length N**)

$$R_N(\Gamma) = \mathbb{E}_{X,Y,x,y} (\hat{y} - y)^2, \text{ where } \dim(Y) = N$$

Pretraining => d^2 -dim linear fitting

$$R_N(\Gamma) = \mathbb{E}_{X,Y,x,y} (\hat{y} - y)^2, \text{ where } \dim(Y) = N \quad \quad \hat{y} = \left\langle \Gamma \cdot \frac{X^\top Y}{N}, x \right\rangle$$

linearly fit $\left(\frac{X^\top Y}{N} \otimes x, y \right)$ with a **matrix parameter** $\Gamma \in \mathbb{R}^{d \times d}$

can we use less than d^2 observations (pretraining tasks)?

Quick facts

$$\hat{y} = \left\langle \Gamma \cdot \frac{X^\top Y}{N}, x \right\rangle$$

$$R_N(\Gamma) = \mathbb{E}_{X,Y,x,y} (\hat{y} - y)^2, \text{ where } \dim(Y) = N$$

- $\Gamma^* = \left(\frac{N+1}{N}H + \frac{\text{tr}(H) + \sigma^2/\psi^2}{N}I \right)^{-1} \approx \left(H + \frac{1}{N}I \right)^{-1}$ $H = \mathbb{E}[xx^\top]$
- $R(\Gamma) - \min R = \left\langle H, (\Gamma - \Gamma^*)\tilde{H}(\Gamma - \Gamma^*)^\top \right\rangle$
- $\tilde{H} = \psi^2 \cdot H \left(\frac{N+1}{N}H + \frac{\text{tr}(H) + \sigma^2/\psi^2}{N}I \right) \approx \psi^2 H \left(H + \frac{1}{N}I \right)$
- $\min R - \sigma^2 = \psi^2 \cdot \text{tr} \left(\Gamma^* H \left(\frac{\text{tr}(H) + \sigma^2/\psi^2}{N} \cdot I + \frac{1}{N}H \right) \right) \lesssim \frac{\psi^2 \text{tr}(H) + \sigma^2}{N}$

$\|\Gamma^*\|_F^2$ could be $\propto N^2d$ if H has many small eigenvalues (e.g., zeros)

Main result 1

Theorem. For T steps of pretraining, we have

$$\mathbb{E}R_N(\Gamma_T) - \min R_N \lesssim \left\langle H\tilde{H}, \left(\prod_{t=1}^T (I - \gamma_t H\tilde{H}) \Gamma^* \right)^2 \right\rangle + (\psi^2 \text{tr}(H) + \sigma^2) \cdot \frac{D_{\text{eff}}}{T_{\text{eff}}}$$

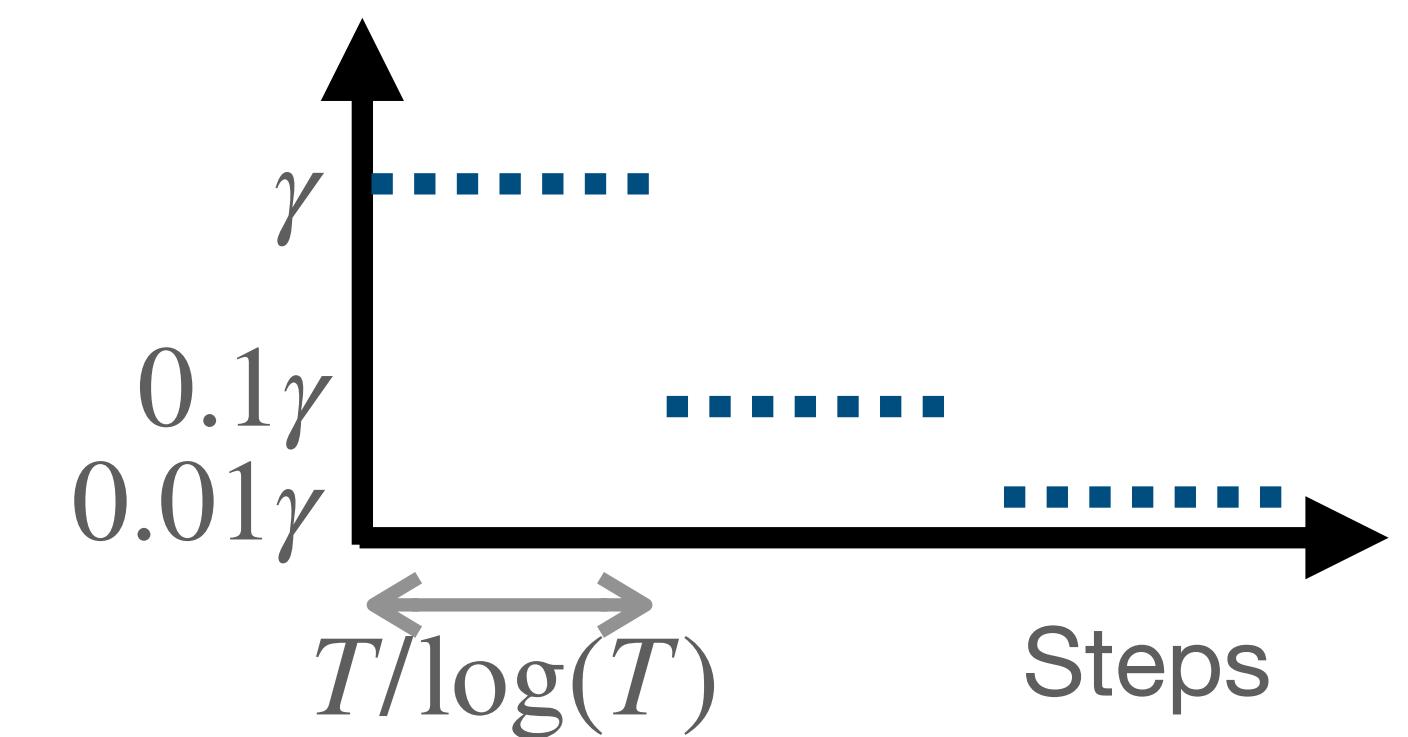
$$T_{\text{eff}} = T/\log(T), \quad D_{\text{eff}} = \sum_{1 \leq i, j \leq d} \min \left\{ 1, \gamma^2 T_{\text{eff}}^2 \lambda_i^2 \tilde{\lambda}_j^2 \right\}$$

λ_i and $\tilde{\lambda}_j$ are eigenvalues of H and \tilde{H}

Corollary. Pretraining enables weak ICL

$$\mathbb{E}R_N(\Gamma_T) = \sigma^2 + \boxed{\min R_N - \sigma^2} + \boxed{\mathbb{E}R_N(\Gamma_T) - \min R_N} \\ \rightarrow 0 \text{ as } N \rightarrow \infty \quad \rightarrow 0 \text{ as } T \rightarrow \infty$$

yay, we match ERM



$$f: \mathbb{R}^{*\times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

Bayesian optimal during inference

Fixing model f , its average risk (conditional on X) at length M is

$$R_M(f; X) = \mathbb{E}_{Y,x,y} (f(X, Y, x) - y)^2, \text{ where } \dim(Y) = M$$

Proposition. Tuned ridge is Bayesian optimal

$$\hat{y} = \langle (X^\top X + \sigma^2/\psi^2 I)^{-1} X^\top Y, x \rangle$$

Moreover, if $\psi^2 \text{tr}(H) \lesssim \sigma^2$, then the average risk (whp) is

$$R_M(\text{ridge}; X) - \sigma^2 \asymp \psi^2 \cdot \sum_i \min\{\lambda_i, \mu_M\}, \text{ where } \mu_M \asymp \frac{\sigma^2/\psi^2}{M}$$

Main result 2

Pretrain at length N

$$\hat{y} = \left\langle \Gamma^* \cdot \frac{X^\top Y}{\dim(Y)}, x \right\rangle \text{ where } \Gamma^* \approx \left(H + \frac{\text{tr}(H) + \sigma^2/\psi^2}{N} I \right)^{-1}$$

Theorem. If $\psi^2 \text{tr}(H) \lesssim \sigma^2$, then the average risk at length M is

$$\mathbb{E}R_M(f; X) - \sigma^2 \approx \boxed{\psi^2 \cdot \sum_i \min\{\lambda_i, \mu_M\}}$$
 ridge risk

$$+ \psi^2 (\mu_M - \mu_N)^2 \cdot \sum_i \min \left\{ \frac{\lambda_i}{\mu_N^2}, \frac{1}{\lambda_i} \right\} \cdot \min \left\{ \frac{\lambda_i}{\mu_M}, 1 \right\}$$

$$\text{where } \mu_M \approx \frac{\sigma^2/\psi^2}{M}, \mu_N \approx \frac{\sigma^2/\psi^2}{N}.$$

We match optimal ridge if $M=N$

$$f: \mathbb{R}^{* \times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$
$$(X, Y, x) \mapsto \hat{y}$$

Four levels of ICL by pretraining

1. “consistent”

ERM, nothing to learn

- error \rightarrow min error when $M, N \rightarrow \infty$

2. “one-point optimal”

one hyperparameter, one thing to learn

- error rate is optimal when $M = N$

supervised learning

3. “uniformly optimal” / algorithm selection

real in-context learning

- error rate is optimal for every $M \leq N$

N hyperparameters
log(N) things to learn

4. “generalizable optimal” / algorithm learning

- error rate is optimal for both $M \leq N$ and $M > N$

a rule to learn

$$f: \mathbb{R}^{*\times d} \otimes \mathbb{R}^* \otimes \mathbb{R}^d \rightarrow \mathbb{R}$$

$$(X, Y, x) \mapsto \hat{y}$$

Where are we

1. “consistent”

- error \rightarrow min error when $M, N \rightarrow \infty$

2. “one-point optimal”

- error rate is optimal when $M = N$



3. “uniformly optimal” / algorithm selection

- error rate is optimal for every $M \leq N$

4. “generalizable optimal” / algorithm learning

- error rate is optimal for both $M \leq N$ and $M > N$

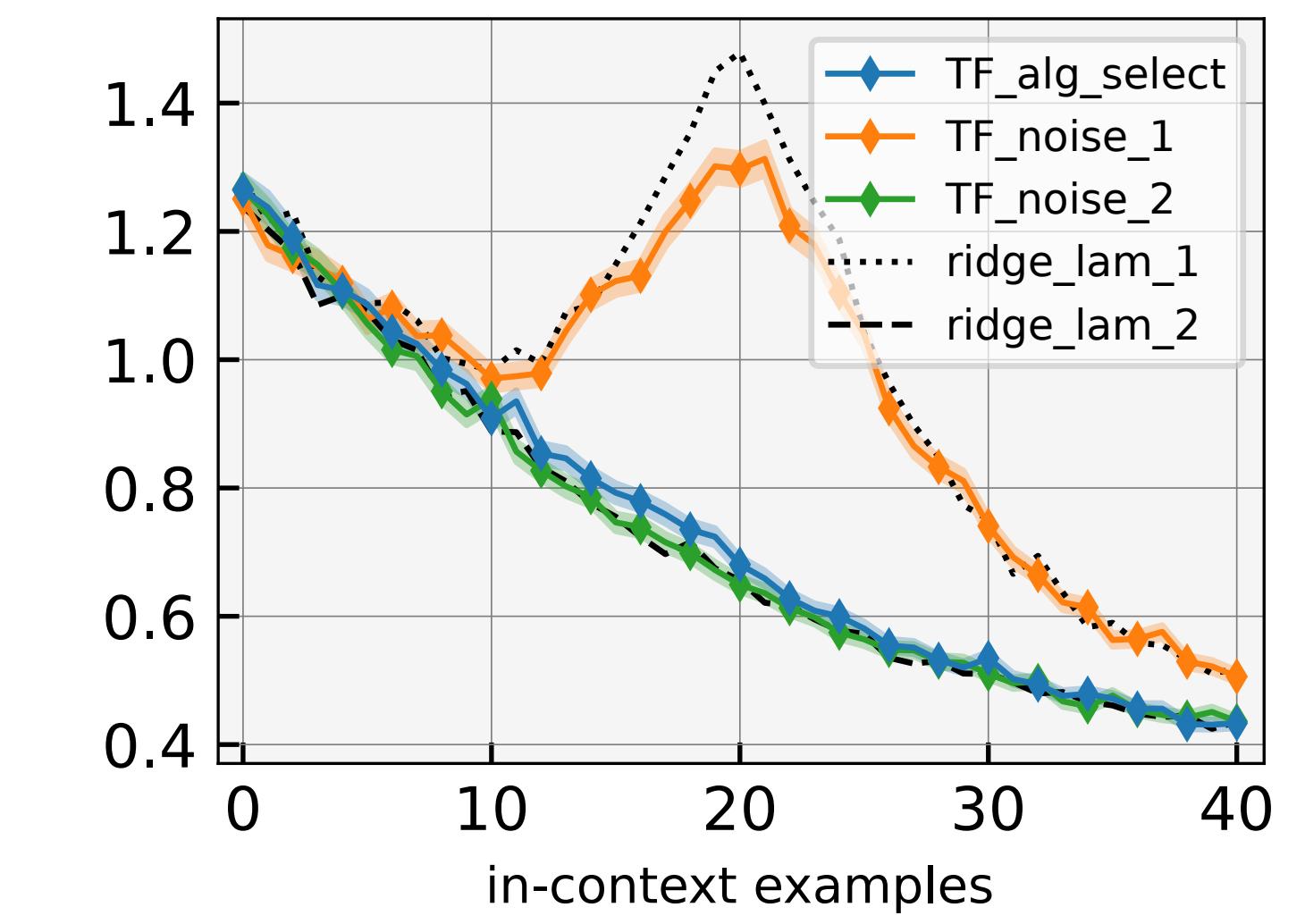
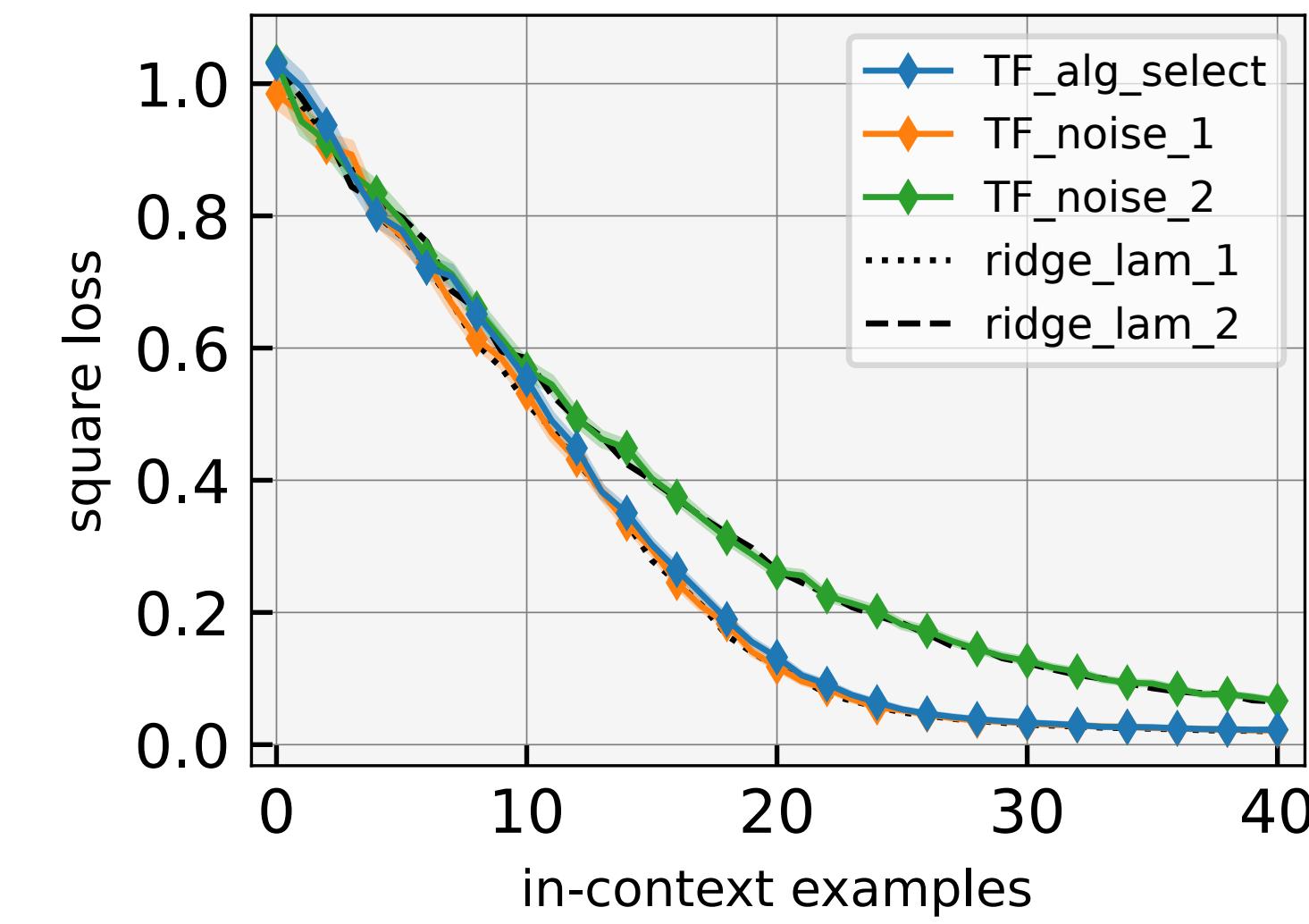
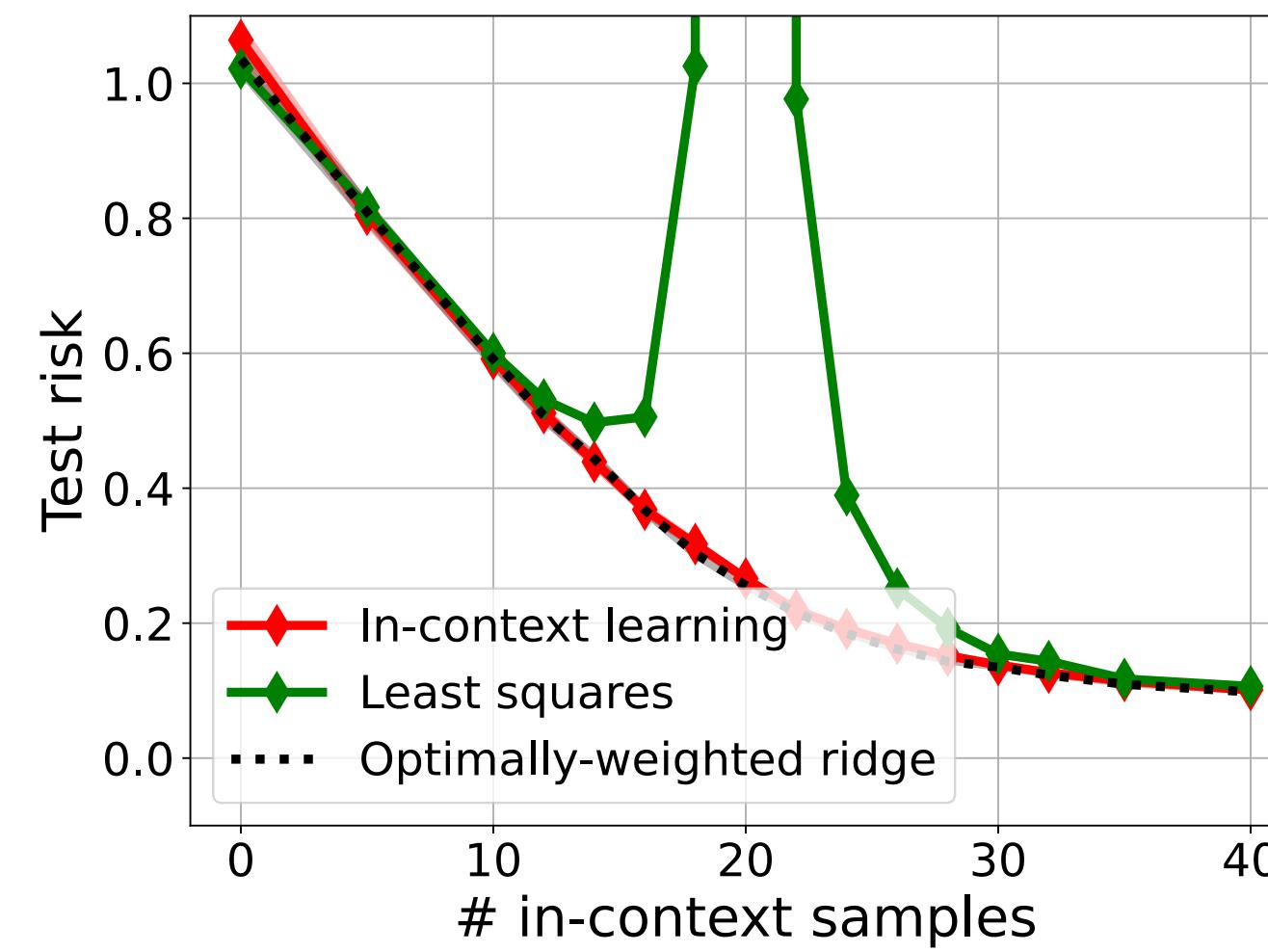
theory \leq level 2

supervised learning

real in-context learning

practice \geq level 3

Evidences for “practice \geq level 3”



Takes

- task complexity of pretraining
- four levels of ICL
- a theory of ICL at level 2

Open questions

- many “obvious” questions
- two less obvious questions
 - can we do ICL theory at level 3?
 - can transformer do ICL at level 4?