The Size of Teachers as a Measure of Data Complexity PAC-Bayes Excess Risk Bounds and Scaling Laws for Neural Networks

Gintare Karolina Dziugaite (Google Deepmind)

Daniel M. Roy (U of Toronto; Vector Institute)

Post-Bayes Online Seminar

We have no good way (yet) of modelling real data, e.g., no notion of the complexity of data.

We have no good way (yet) of modelling real data, e.g., no notion of the complexity of data.

OTOH, empirical "neural scaling laws" predict performance across wide ranges of data and model sizes.

We have no good way (yet) of modelling real data, e.g., no notion of the complexity of data.

OTOH, empirical "neural scaling laws" predict performance across wide ranges of data and model sizes.

What might neural scaling laws tell us about the complexity of the underlying data?



Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C_{\mu}(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\mathsf{exp}(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1} n). \end{split}$$

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C_{\mu}(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1}n). \end{split}$$

Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then...

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C_{\mu}(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1}n). \end{split}$$

Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then... $C_{\mu}(\varepsilon) \not\in O(\varepsilon^{-p})$.

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C_{\mu}(\varepsilon) = O(\text{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\exp(\text{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1} n). \end{split}$$

Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then... $C_{\mu}(\varepsilon) \not\in O(\varepsilon^{-p})$.

1. Such function $C_{\mu}(\cdot)$ can be viewed as defining a notion of complexity of the data distribution $\mu.$

Let $r(\theta)$ be the risk (say, probability of misclassification) of a predictor θ .

Let $\hat{\theta}_n$ be the predictor produced by our learning algorithm given n training data training a model of size m(n).

Theorem. For every data distribution μ , there exists a function $C_{\mu}(\varepsilon)$ such that, with high probability under $\mu^{\otimes n}$,

$$r(\hat{\theta}_n) \leq \inf_{\varepsilon} \left\{ \varepsilon + O\bigg(\sqrt{\frac{\varepsilon \cdot C_{\mu}(\varepsilon)}{n}} \ \bigg) \right\}.$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C_{\mu}(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C_{\mu}(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1}n). \end{split}$$

Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then... $C_{\mu}(\varepsilon) \not\in O(\varepsilon^{-p})$.

- 1. Such function $C_{\mu}(\cdot)$ can be viewed as defining a notion of complexity of the data distribution $\mu.$
 - 2. Scaling laws can provide evidence of complexity, in view of upper bounds.

We propose to measure the complexity of data in terms of ... the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ . Can be viewed as a rate–distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

▶ Study toy model of neural network training: a randomly initialized neural networks that fit the data.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

Study toy model of neural network training: a randomly initialized neural networks that fit the data.

Assume that some "teacher" network has zero risk. ("realizable" setting)

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
- **Caveat:** We'll quantize the weights (or need to introduce some notion of margin.)

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ▶ Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - **Caveat:** We'll quantize the weights (or need to introduce some notion of margin.)

Primi

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ▶ Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ▶ Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
- Caveat: We'll quantize the weights (or need to introduce some notion of margin.)

Primi

- ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
- Drop assumption that some teacher network has zero risk. ("agnostic" setting)
- ► Conclusion: Number of teacher parameters determines sample complexity (for excess risk).

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- ▶ Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
- Caveat: We'll quantize the weights (or need to introduce some notion of margin.)

Primi

- ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
- Drop assumption that some teacher network has zero risk. ("agnostic" setting)
- ► Conclusion: Number of teacher parameters determines sample complexity (for excess risk).

Bonus: Nonvacuous bounds for MNIST.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ▶ Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - ► Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi
 - Turn agnostic bound into an oracle inequality: Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ► Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - ► Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi
 - Turn agnostic bound into an oracle inequality: Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.
 - $\qquad \qquad \textbf{Introduce } C(\varepsilon) \text{ and rewrite bound, obtain } \inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \} \text{ bound.}$

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
- Caveat: We'll quantize the weights (or need to introduce some notion of margin.)

Primi

- ► Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
- Drop assumption that some teacher network has zero risk. ("agnostic" setting)
- Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
- Bonus: Nonvacuous bounds for MNIST.

Secondi

- Turn agnostic bound into an oracle inequality:
 Risk of Gibbs posterior sample no more than any teacher's ri
- Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.
- $\qquad \qquad \textbf{Introduce } C(\varepsilon) \text{ and rewrite bound, obtain } \inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \} \text{ bound.}$
- ▶ Conclusion: Scaling laws suggest $C(\varepsilon) \in \omega(\varepsilon^{-p})$ for some p.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - ► Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ► Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi
 - Turn agnostic bound into an oracle inequality: Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.
 - ▶ Introduce $C(\varepsilon)$ and rewrite bound, obtain $\inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \}$ bound.
 - ▶ Conclusion: Scaling laws suggest $C(\varepsilon) \in \omega(\varepsilon^{-p})$ for some p.
- Dolce

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ► Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - ► Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi
 - Turn agnostic bound into an oracle inequality: Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.
 - Introduce $C(\varepsilon)$ and rewrite bound, obtain $\inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \}$ bound.
 - ▶ Conclusion: Scaling laws suggest $C(\varepsilon) \in \omega(\varepsilon^{-p})$ for some p.
- Dolce
 - Don't want teachers but distributions on teachers.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

Antipasti

- Study toy model of neural network training: a randomly initialized neural networks that fit the data.
- Assume that some "teacher" network has zero risk. ("realizable" setting)
- ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
- ► Caveat: We'll quantize the weights (or need to introduce some notion of margin.)

Primi

- ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
- Drop assumption that some teacher network has zero risk. ("agnostic" setting)
- Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
- Bonus: Nonvacuous bounds for MNIST.

Secondi

- ► Turn agnostic bound into an oracle inequality:
 - Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.
- ▶ Introduce $C(\varepsilon)$ and rewrite bound, obtain $\inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \}$ bound.
- ▶ Conclusion: Scaling laws suggest $C(\varepsilon) \in \omega(\varepsilon^{-p})$ for some p.
- Dolce
 - Don't want teachers but distributions on teachers.
 - Distributions avoid quantization/discretization. (Bits back encoding.) BONUS: New perspective on variational inference.

We propose to measure the complexity of data in terms of ...

the size $C_{\mu}(\varepsilon)$ of the smallest network that achieves ε risk under μ .

Can be viewed as a rate-distortion function, specific to the data distribution (Hafez-Kolahi, Moniri, Kasaei 2024; Hafez-Kolahi, Moniri et al. 2021)

- Antipasti
 - Study toy model of neural network training: a randomly initialized neural networks that fit the data.
 - Assume that some "teacher" network has zero risk. ("realizable" setting)
 - ► Conclusion (Buzaglo et al. '24): Number of teacher parameters determines sample complexity.
 - ► Caveat: We'll quantize the weights (or need to introduce some notion of margin.)
- Primi
 - ▶ Study less toy model: sample from the Gibbs posterior $\propto \exp\{-\beta \hat{r}_n(\theta) + d\pi(\theta)\}$.
 - Drop assumption that some teacher network has zero risk. ("agnostic" setting)
 - Conclusion: Number of teacher parameters determines sample complexity (for excess risk).
 - Bonus: Nonvacuous bounds for MNIST.
- Secondi
 - Turn agnostic bound into an oracle inequality:

Risk of Gibbs posterior sample no more than any teacher's risk plus a size penalty.

- Introduce $C(\varepsilon)$ and rewrite bound, obtain $\inf_{\varepsilon} \{ \varepsilon + C(\varepsilon) ... \}$ bound.
- ▶ Conclusion: Scaling laws suggest $C(\varepsilon) \in \omega(\varepsilon^{-p})$ for some p.
- Dolce
 - Don't want teachers but distributions on teachers.
 - Distributions avoid quantization/discretization. (Bits back encoding.) BONUS: New perspective on variational inference.
 - **Conclusion:** $C(\varepsilon)$ more complicated.

Learning Setup

Labelled data: denoted by z, z_i, Z, \ldots

Predictors: identified with parameters, denoted by $\theta, \hat{\theta}, \dots$

Loss: $\ell(\theta, z)$, e.g., $\ell(\theta, (x, y)) = \mathbb{I}(\text{network with weights } \theta \text{ on input } x \text{ does not output label } y)$

Data distribution: μ , presumed unknown

Risk: $r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$

Data: $S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$

Empirical Risk: $\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$

Prior: distribution π on Θ , which is nonrandom.

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Buzaglo et al.'s assumptions

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk).

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Fix a probability measure π on Θ .

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Fix a probability measure π on Θ . Let $\theta \sim \pi$. If $\theta \in \hat{\Theta}_0$, accept and set $\hat{\theta} = \theta$. Otherwise, try again.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Fix a probability measure π on Θ . Let $\theta \sim \pi$. If $\theta \in \hat{\Theta}_0$, accept and set $\hat{\theta} = \theta$. Otherwise, try again.

Distribution of $\hat{\theta}$ is the "posterior" $\hat{\rho} = \pi(\cdot \mid \hat{\Theta}_0)$ given by $\hat{\rho}(A) = \frac{\pi(A \cap \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Fix a probability measure π on Θ . Let $\theta \sim \pi$. If $\theta \in \hat{\Theta}_0$, accept and set $\hat{\theta} = \theta$. Otherwise, try again.

Distribution of $\hat{\theta}$ is the "posterior" $\hat{\rho} = \pi(\,\cdot\mid\hat{\Theta}_0)$ given by $\hat{\rho}(A) = \frac{\pi(A\cap\hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$. Equivalently, $\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta\in\hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$.

Buzaglo et al.'s assumptions

Let
$$\Theta^* = \{\theta \in \Theta : r(\theta^*) = 0\}.$$

Assumption (realizability). There exists $\theta^* \in \Theta^*$. Assumption (finiteness). Θ is a finite set.

Buzaglo et al.: What's the risk of a "random interpolator"?

Let $\hat{\Theta}_0 = \{\theta \in \Theta : \hat{r}_n(\theta) = 0\}$ be the set of interpolating predictors (i.e., with zero empirical risk). Note: $\hat{\Theta}_0 \supseteq \Theta^*$.

We want a random element from $\hat{\Theta}_0$. How?

Fix a probability measure π on Θ . Let $\theta \sim \pi$. If $\theta \in \hat{\Theta}_0$, accept and set $\hat{\theta} = \theta$. Otherwise, try again.

Distribution of $\hat{\theta}$ is the "posterior" $\hat{\rho} = \pi(\,\cdot\mid\hat{\Theta}_0)$ given by $\hat{\rho}(A) = \frac{\pi(A\cap\hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$. Equivalently, $\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta\in\hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$.

So... what's the risk of $\hat{\theta}$ sampled from the posterior $\hat{\rho}$?

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Let $\hat{\theta}$ be a sample from $\hat{\rho}$. (That is, $\hat{\theta} \mid S \sim \hat{\rho}$.)

$${\bf Risk:} \quad r(\theta) = \int \ell(\theta,z) \mu({\rm d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Let $\hat{\theta}$ be a sample from $\hat{\rho}$. (That is, $\hat{\theta} \mid S \sim \hat{\rho}$.)

PAC-Bayes offers comparisons for the random variables

$$r(\hat{\theta})$$
 and $\hat{r}_n(\hat{\theta})$ and $\underbrace{\log \frac{\mathrm{d} \rho}{\mathrm{d} \pi}(\hat{\theta})}_{\text{information density}}$

$$\textbf{Risk:} \quad r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Let $\hat{\theta}$ be a sample from $\hat{\rho}$. (That is, $\hat{\theta} \mid S \sim \hat{\rho}$.)

PAC-Bayes offers comparisons for the random variables

$$r(\hat{\theta})$$
 and $\hat{r}_n(\hat{\theta})$ and $\log \frac{d\rho}{d\pi}(\hat{\theta})$

Classical case controls expectations of these quantities under $\hat{\rho}$.

7/29

$$\textbf{Risk:} \quad r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Let $\hat{\theta}$ be a sample from $\hat{\rho}$. (That is, $\hat{\theta} \mid S \sim \hat{\rho}$.)

PAC-Bayes offers comparisons for the random variables

$$r(\hat{\theta})$$
 and $\hat{r}_n(\hat{\theta})$ and $\log \frac{d\rho}{d\pi}(\hat{\theta})$

Classical case controls expectations of these quantities under $\hat{\rho}$.

7/29

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Risk:
$$r(\theta) = \int \ell(\theta, z) \mu(dz)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Theorem (Single sample bound; Catoni 2007)

Fix $\lambda > 0$. With probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \bigg(\hat{r}_n(\hat{\theta}) + \frac{\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\hat{\theta}) + \log \frac{1}{\delta}}{\lambda} \bigg) \qquad \textit{where } \Phi_a^{-1}(q) = (1 - \exp(-aq))/(1 - \exp(-a)).$$

Risk:
$$r(\theta) = \int \ell(\theta, z) \mu(dz)$$

Data:
$$S = (Z_1, \ldots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Theorem (Single sample bound; Catoni 2007)

Fix $\lambda > 0$. With probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \bigg(\hat{r}_n(\hat{\theta}) + \frac{\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\hat{\theta}) + \log \frac{1}{\delta}}{\lambda} \bigg) \qquad \textit{where } \Phi_a^{-1}(q) = (1 - \exp(-aq))/(1 - \exp(-a)).$$

$$\text{How should we interpret } \Phi_{\lambda/n}^{-1} \textbf{?} \qquad \inf_{\lambda \in \mathbb{R}_+} \Phi_{\lambda/n}^{-1} \Big(q + \frac{d}{\lambda} \Big) \leq q + \frac{2d}{n} + \sqrt{\frac{2dq}{n}}, \text{ provided r.h.s. is less than } 1/2.$$

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Corollary

Assume $\hat{r}_n(\hat{\theta}) = 0$ a.s. under $\hat{\theta} \mid S \sim \hat{\rho}$.

Risk:
$$r(\theta) = \int \ell(\theta,z) \mu(\mathrm{d}z)$$

Data:
$$S = (Z_1, \dots, Z_n) \sim \mu^{\otimes n}$$

Empirical Risk:
$$\hat{r}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{\theta}, z_i)$$

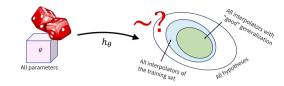
Posterior: distribution $\hat{\rho}$ on Θ , which may depend on S.

Prior: distribution π on Θ , which is nonrandom.

Corollary

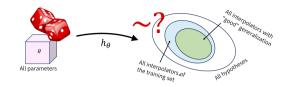
Assume $\hat{r}_n(\hat{\theta}) = 0$ a.s. under $\hat{\theta} \mid S \sim \hat{\rho}$. With probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \frac{\log \frac{d\hat{\rho}}{d\pi}(\hat{\theta}) + \log \frac{1}{\delta}}{n}$$



Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta \in \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$



Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta \in \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$

By construction, $\hat{r}_n(\theta)=0$ and $\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta)=\log \frac{1}{\pi(\Theta^*)}$ for $\hat{\rho}$ -almost all θ .

Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$rac{\mathsf{d}\hat{
ho}}{\mathsf{d}\pi}(heta) = rac{\mathbb{I}(heta\in\hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$

By construction, $\hat{r}_n(\theta)=0$ and $\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta)=\log \frac{1}{\pi(\Theta^*)}$ for $\hat{\rho}$ -almost all θ .

What's the risk of $\hat{\theta}$ sampled from $\hat{\rho}$?

Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta \in \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$

By construction, $\hat{r}_n(\theta)=0$ and $\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta)=\log \frac{1}{\pi(\Theta^*)}$ for $\hat{\rho}$ -almost all θ .

What's the risk of $\hat{\theta}$ sampled from $\hat{\rho}$?

Lemma (Risk of posterior sampling; Dziugaite & R. 2025)

Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{\mathbb{I}(\theta \in \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$

By construction, $\hat{r}_n(\theta)=0$ and $\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta)=\log \frac{1}{\pi(\Theta^*)}$ for $\hat{\rho}$ -almost all θ .

What's the risk of $\hat{\theta}$ sampled from $\hat{\rho}$?

Lemma (Risk of posterior sampling; Dziugaite & R. 2025)

Fix π and assume π -realizability. Let $\hat{\theta} \mid S \sim \hat{\rho}$. With probability at least $1 - \delta$,

Recall: $\hat{\rho}$ is posterior given interpolation, given by

$$\frac{\mathsf{d}\hat{\rho}}{\mathsf{d}\pi}(\theta) = \frac{\mathbb{I}(\theta \in \hat{\Theta}_0)}{\pi(\hat{\Theta}_0)}$$

By construction, $\hat{r}_n(\theta)=0$ and $\log \frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta)=\log \frac{1}{\pi(\Theta^*)}$ for $\hat{\rho}$ -almost all θ .

What's the risk of $\hat{\theta}$ sampled from $\hat{\rho}$?

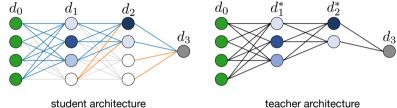
Lemma (Risk of posterior sampling; Dziugaite & R. 2025)

Fix π and assume π -realizability. Let $\hat{\theta} \mid S \sim \hat{\rho}$. With probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \frac{\log \frac{1}{\pi(\Theta^*)} + \log \frac{1}{\delta}}{n}$$

Antipasti: Buzaglo et al.'s lower bound on $\pi(\Theta^*)$, the probability of interpolating

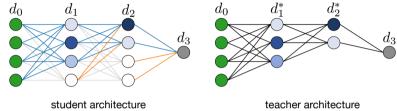
For a network θ , let $m(\theta)$ be the smallest number of parameters we must specify to produce a network functionally equivalent to θ .



Note: $m(\theta) = \#$ parameters in teacher network + # neurons in student network

Antipasti: Buzaglo et al.'s lower bound on $\pi(\Theta^*)$, the probability of interpolating

For a network θ , let $m(\theta)$ be the smallest number of parameters we must specify to produce a network functionally equivalent to θ .

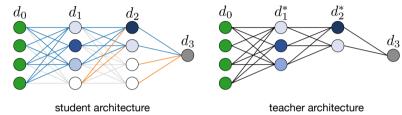


Note: $m(\theta) = \#$ parameters in teacher network + # neurons in student network

Assumption (quantization). Weights of all networks are quantized into Q levels, and zero is one level.

Antipasti: Buzaglo et al.'s lower bound on $\pi(\Theta^*)$, the probability of interpolating

For a network θ , let $m(\theta)$ be the smallest number of parameters we must specify to produce a network functionally equivalent to θ .

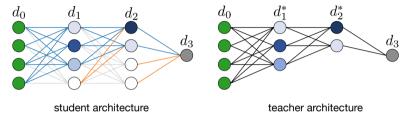


Note: $m(\theta) = \#$ parameters in teacher network + # neurons in student network

Assumption (quantization). Weights of all networks are quantized into Q levels, and zero is one level. **Assumption (batch-norm-like init).** Weights leaving each neuron are multiplied by a neuron-specific scaling weight.

Antipasti: Buzaglo et al.'s lower bound on $\pi(\Theta^*)$, the probability of interpolating

For a network θ , let $m(\theta)$ be the smallest number of parameters we must specify to produce a network functionally equivalent to θ .



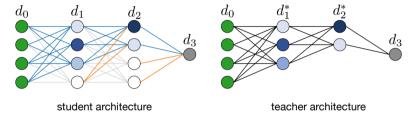
Note: $m(\theta) = \#$ parameters in teacher network + # neurons in student network

Assumption (quantization). Weights of all networks are quantized into Q levels, and zero is one level. **Assumption (batch-norm-like init).** Weights leaving each neuron are multiplied by a neuron-specific scaling weight. **Assumption (uniform prior).** π is uniform distribution over all quantized student networks.

Antipasti: Buzaglo et al.'s lower bound on $\pi(\Theta^*)$, the probability of interpolating

For a network θ , let $m(\theta)$ be the smallest number of parameters we must specify to produce a network functionally

equivalent to θ .



Note: $m(\theta) = \#$ parameters in teacher network + # neurons in student network

Assumption (quantization). Weights of all networks are quantized into Q levels, and zero is one level. **Assumption (batch-norm-like init).** Weights leaving each neuron are multiplied by a neuron-specific scaling weight. **Assumption (uniform prior).** π is uniform distribution over all quantized student networks.

Second contribution of Buzaglo et al.

Let
$$\theta^* \in \Theta^*$$
. Then $\log \frac{1}{\pi(\Theta^*)} = O(m(\theta^*) \log Q)$.

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Assume we have quantization, batch-norm-like init, a uniform prior π , and π -realizability.

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Assume we have quantization, batch-norm-like init, a uniform prior π , and π -realizability.

Let $m^* = \inf_{\theta^* \in \Theta^*} m(\theta^*)$ be the smallest teacher network in Θ^* .

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Assume we have quantization, batch-norm-like init, a uniform prior π , and π -realizability.

Let $m^* = \inf_{\theta^* \in \Theta^*} m(\theta^*)$ be the smallest teacher network in Θ^* .

Let S be n i.i.d. data and let $\hat{\theta}$ be a random interpolating network.

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Assume we have quantization, batch-norm-like init, a uniform prior π , and π -realizability.

Let $m^* = \inf_{\theta^* \in \Theta^*} m(\theta^*)$ be the smallest teacher network in Θ^* .

Let S be n i.i.d. data and let $\hat{\theta}$ be a random interpolating network. With probability at least $1-\delta$,

Theorem (Buzaglo et al. 2024; Dziugaite & R. 2025)

Assume we have quantization, batch-norm-like init, a uniform prior π , and π -realizability.

Let $m^* = \inf_{\theta^* \in \Theta^*} m(\theta^*)$ be the smallest teacher network in Θ^* .

Let S be n i.i.d. data and let $\hat{\theta}$ be a random interpolating network.

With probability at least $1 - \delta$,

$$r(\hat{\theta}) \le \frac{m^* \log Q + \log \frac{1}{\delta}}{n}$$

Equivalently, if we have at least

$$n \geq \frac{m^* \log Q + \log \frac{1}{\delta}}{\varepsilon}$$

samples, then $r(\hat{\theta}) \leq \varepsilon$ with probability at least $1 - \delta$.



Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

• Student and teacher have the same depth.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.
- Realizability assumption not met in practice.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.
- Realizability assumption not met in practice.
- Quantization assumption doesn't match practice, and is essential to the argument.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.
- Realizability assumption not met in practice.
- Quantization assumption doesn't match practice, and is essential to the argument.
- Predicts faster rate than seen empirically (in scaling laws).

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.
- Realizability assumption not met in practice.
- Quantization assumption doesn't match practice, and is essential to the argument.
- Predicts faster rate than seen empirically (in scaling laws).
- Cannot explain performance as networks diverge in size.

Nice observation that parametrization can induce "implicit bias", even from "uniform" sampling.

On the other hand...

- Student and teacher have the same depth.
- The parametrization we exploited is not necessary for generalization.
- Posterior sampling may not be a good model.
- Realizability assumption not met in practice.
- Quantization assumption doesn't match practice, and is essential to the argument.
- Predicts faster rate than seen empirically (in scaling laws).
- Cannot explain performance as networks diverge in size.

We'll start with realizability (which will see us rethink posterior sampling, rates, etc.).

Non-realizable case means Θ^* empty.

Non-realizable case means Θ^* empty. Then possibly $\pi(\hat{\Theta}_0)=0$ with positive probability, hence...

Non-realizable case means Θ^* empty. Then possibly $\pi(\hat{\Theta}_0) = 0$ with positive probability, hence... no posterior.

Non-realizable case means Θ^* empty. Then possibly $\pi(\hat{\Theta}_0)=0$ with positive probability, hence... no posterior.

Gibbs posteriors allow us to model soft (stochastic) optimization and is well-defined in the non-realizable case.

Non-realizable case means Θ^* empty. Then possibly $\pi(\hat{\Theta}_0) = 0$ with positive probability, hence... no posterior.

Gibbs posteriors allow us to model soft (stochastic) optimization and is well-defined in the non-realizable case.

Defn. The *Gibbs posterior* for a prior π on Θ , inverse temperature $\beta > 0$, and empirical risk \hat{r}_n under data S, is the distribution $\hat{\rho}$ on Θ dominated by π given by

$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\} \propto \exp\{-\beta\,\hat{r}_n(\theta)\},$$

where $Z^\pi_\beta(S) := \int \exp\{-\beta \, \hat{r}_n(\theta)\} \pi(\mathrm{d}\theta)$ is the normalization constant.

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

• Generalizes hard posterior. For fixed S, converges to (hard) posterior as inverse temperature $\beta \to \infty$.

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

- Generalizes hard posterior. For fixed S, converges to (hard) posterior as inverse temperature $\beta \to \infty$.
- Related to gradient flow + noise. Assuming π is absolutely continuous with density p, the Gibbs posterior is the stationary distribution of Langevin dynamics with drift $\beta \hat{r}_n(\theta) + \log p(\theta)$.

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

- Generalizes hard posterior. For fixed S, converges to (hard) posterior as inverse temperature $\beta \to \infty$.
- Related to gradient flow + noise. Assuming π is absolutely continuous with density p, the Gibbs posterior is the stationary distribution of Langevin dynamics with drift $\beta \hat{r}_n(\theta) + \log p(\theta)$.
- Related to SGD + noise + decaying stepsize. Gibbs posterior is also the limiting dynamics for Stochastic Gradient Langevin Dynamics (SGD + noise) for decaying step sizes.

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

- Generalizes hard posterior. For fixed S, converges to (hard) posterior as inverse temperature $\beta \to \infty$.
- Related to gradient flow + noise. Assuming π is absolutely continuous with density p, the Gibbs posterior is the stationary distribution of Langevin dynamics with drift $\beta \hat{r}_n(\theta) + \log p(\theta)$.
- Related to SGD + noise + decaying stepsize. Gibbs posterior is also the limiting dynamics for Stochastic Gradient Langevin Dynamics (SGD + noise) for decaying step sizes.
- Related to SGD + noise + fixed step size in large data limit. The limiting OU process dynamics for SGLD and SGD
 in fixed step size regimes, where data, # of total iterations diverges. (These results don't cover high-dimensional case
 yet, though.)

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta).$

- Generalizes hard posterior. For fixed S, converges to (hard) posterior as inverse temperature $\beta \to \infty$.
- Related to gradient flow + noise. Assuming π is absolutely continuous with density p, the Gibbs posterior is the stationary distribution of Langevin dynamics with drift $\beta \hat{r}_n(\theta) + \log p(\theta)$.
- Related to SGD + noise + decaying stepsize. Gibbs posterior is also the limiting dynamics for Stochastic Gradient Langevin Dynamics (SGD + noise) for decaying step sizes.
- Related to SGD + noise + fixed step size in large data limit. The limiting OU process dynamics for SGLD and SGD
 in fixed step size regimes, where data, # of total iterations diverges. (These results don't cover high-dimensional case
 yet, though.)
- Gibbs posteriors optimize PAC-Bayes bounds. The Gibbs posteriors above arises as the solution to variational problems: $\arg\min_{\rho} \rho(\hat{r}_n) + \beta^{-1} KL(\rho|\pi)$. That is, it minimizes certain PAC-Bayes bounds.

Primi: Single-sample PAC-Bayes Bound for Gibbs Posterior

Gibbs posterior:
$$\frac{\mathrm{d}\hat{\rho}}{\mathrm{d}\pi}(\theta) = \frac{1}{Z_{\beta}^{\pi}(S)} \exp\{-\beta\,\hat{r}_n(\theta)\}$$
 where $Z_{\beta}^{\pi}(S) := \int \exp\{-\beta\,\hat{r}_n(\theta)\}\pi(\mathrm{d}\theta) = \pi(\exp\{-\beta\hat{r}_n\}).$

Corollary

Let $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature β , and \hat{r}_n , and let $\hat{\theta} \mid S \sim \hat{\rho}$. With probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[(1-\lambda^{-1}\beta) \hat{r}_n(\hat{\theta}) - \lambda^{-1} \log \left(Z_\beta^\pi(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

If
$$\beta \geq \lambda$$
,

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\beta}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

Primi: Back to Teachers

Taking inverse temperature $\beta = \lambda$, we have a risk bound

Primi: Back to Teachers

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

We will see that this term is the **analogue of** $\lambda^{-1}\log\frac{1}{\pi(\Theta^*)}.$

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

We will see that this term is the analogue of $\lambda^{-1}\log\frac{1}{\pi(\Theta^*)}.$

Teacher lower bound on local entropy

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

We will see that this term is the analogue of $\lambda^{-1}\log\frac{1}{\pi(\Theta^*)}.$

Teacher lower bound on local entropy

Let $\theta^* \in \Theta$ be an (arbitrary!) teacher. Let $E_{\theta^*} = \{\theta \in \Theta : \theta^* \text{ and } \theta \text{ have same minimal model}\}.$

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \left[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \right].$$

We will see that this term is the analogue of $\lambda^{-1}\log\frac{1}{\pi(\Theta^*)}.$

Teacher lower bound on local entropy

Let $\theta^* \in \Theta$ be an (arbitrary!) teacher. Let $E_{\theta^*} = \{\theta \in \Theta : \theta^* \text{ and } \theta \text{ have same minimal model}\}$.

$$Z_{\lambda}^{\pi}(S) = \int \exp\{-\lambda \, \hat{r}_n(\theta)\} \, \pi(\mathsf{d}\theta) \ge \exp\{-\lambda \, \hat{r}_n(\theta^*)\} \, \pi(E_{\theta^*}). \tag{1}$$

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \leq \Phi_{\lambda/n}^{-1} \Big[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \Big].$$

We will see that this term is the **analogue of** $\lambda^{-1}\log\frac{1}{\pi(\Theta^*)}$.

Teacher lower bound on local entropy

Let $\theta^* \in \Theta$ be an (arbitrary!) teacher. Let $E_{\theta^*} = \{\theta \in \Theta : \theta^* \text{ and } \theta \text{ have same minimal model}\}$.

$$Z_{\lambda}^{\pi}(S) = \int \exp\{-\lambda \, \hat{r}_n(\theta)\} \, \pi(\mathsf{d}\theta) \ge \exp\{-\lambda \, \hat{r}_n(\theta^*)\} \, \pi(E_{\theta^*}). \tag{1}$$

Thus, the local entropy obeys

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) \le -\lambda^{-1}\log\left(\exp\{-\lambda\,\hat{r}_n(\theta^*)\}\,\pi(E_{\theta^*})\right) \tag{2}$$

$$= \hat{r}_n(\theta^*) + \lambda^{-1} \log \frac{1}{\pi(E_{\theta^*})}$$
 (3)

Taking inverse temperature $\beta = \lambda$, we have a risk bound

$$r(\hat{\theta}) \le \Phi_{\lambda/n}^{-1} \left[-\lambda^{-1} \log \left(Z_{\lambda}^{\pi}(S) \right) + \lambda^{-1} \log \frac{1}{\delta} \right].$$

We will see that this term is the **analogue of** $\lambda^{-1} \log \frac{1}{\pi(\Theta^*)}$.

Teacher lower bound on local entropy

Let $\theta^* \in \Theta$ be an (arbitrary!) teacher. Let $E_{\theta^*} = \{\theta \in \Theta : \theta^* \text{ and } \theta \text{ have same minimal model}\}$.

$$Z_{\lambda}^{\pi}(S) = \int \exp\{-\lambda \, \hat{r}_n(\theta)\} \, \pi(\mathrm{d}\theta) \ge \exp\{-\lambda \, \hat{r}_n(\theta^*)\} \, \pi(E_{\theta^*}). \tag{1}$$

Thus, the local entropy obeys

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) \leq -\lambda^{-1}\log\left(\exp\{-\lambda\,\hat{r}_{n}(\theta^{*})\}\,\pi(E_{\theta^{*}})\right)$$

$$= \hat{r}_n(\theta^*) + \lambda^{-1} \log \frac{1}{\pi(F_{n-1})}$$

(2)

(3)

Indeed, the bound holds simultaneously for all
$$\theta^* \in \Theta^*$$
, and so we have

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) \leq \inf_{\theta^{*}}\left\{\hat{r}_{n}(\theta^{*}) + \lambda^{-1}\log\frac{1}{\pi(E_{\theta^{*}})}\right\} \tag{4}$$

Key term is again $\log \frac{1}{\pi(E_{\theta^*})}$

Key term is again $\log \frac{1}{\pi(E_{\theta^*})} \leq m(\theta^*) \log Q$, by the same argument as Buzaglo et al.

Key term is again $\log \frac{1}{\pi(E_{\theta^*})} \leq m(\theta^*) \log Q$, by the same argument as Buzaglo et al.

$$\text{Take } \theta^{**} \in \arg \min_{\theta \in \Theta^*} \Big\{ r(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \Big\}.$$

Key term is again $\log \frac{1}{\pi(E_{0*})} \leq m(\theta^*) \log Q$, by the same argument as Buzaglo et al.

$$\mathrm{Take}\ \theta^{**} \in \mathrm{arg} \, \mathrm{min}_{\theta \in \Theta^*} \, \Big\{ r(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \Big\}.$$

With probability at least $1 - \delta$,

$$\begin{split} \inf_{\theta^* \in \Theta} \left\{ \hat{r}_n(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \right\} &\leq r(\theta^{**}) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^{**}) \log Q \\ &\leq \inf_{\theta^* \in \Theta} \left\{ r(\theta^*) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^*) \log Q \right\} \end{split}$$

Key term is again $\log \frac{1}{\pi(E_{a*})} \leq m(\theta^*) \log Q$, by the same argument as Buzaglo et al.

 $\operatorname{Take} \theta^{**} \in \arg \min_{\theta \in \Theta^*} \Big\{ r(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \Big\}.$

With probability at least $1 - \delta$,

$$\begin{split} \inf_{\theta^* \in \Theta} \left\{ \hat{r}_n(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \right\} &\leq r(\theta^{**}) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^{**}) \log Q \\ &\leq \inf_{\theta^* \in \Theta} \left\{ r(\theta^*) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^*) \log Q \right\} \end{split}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

Key term is again $\log \frac{1}{\pi(E_{a*})} \leq m(\theta^*) \log Q$, by the same argument as Buzaglo et al.

 $\operatorname{Take} \theta^{**} \in \operatorname{arg\,min}_{\theta \in \Theta^*} \Big\{ r(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \Big\}.$

With probability at least $1 - \delta$,

$$\begin{split} \inf_{\theta^* \in \Theta} \left\{ \hat{r}_n(\theta^*) + \lambda^{-1} m(\theta^*) \log Q \right\} &\leq r(\theta^{**}) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^{**}) \log Q \\ &\leq \inf_{\theta^* \in \Theta} \left\{ r(\theta^*) + O(\sqrt{n^{-1} \log 1/\delta}) + \lambda^{-1} m(\theta^*) \log Q \right\} \end{split}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\theta^* \in \Theta} O\bigg(r(\theta^*) + \sqrt{\frac{m(\theta^*)\log Q + \log 1/\delta}{n}}\,\bigg)$$

Can get a bound of the form $r + \frac{c}{n} + \sqrt{\frac{rc}{n}}$ too.

Primi: Toy MNIST Experiment

We validate our bounds on the MNIST dataset.

Experimental setup:

- 2-layer MLP with 3 hidden units (3167 parameters)
- Trained on MNIST dataset
- Obtained θ^* with 0.145 risk
- Fast rate bound yields 0.374 risk had we been able to quantize to 4-bits... but we haven't tried.

Primi: Takeaways

This bound balances teacher risk and complexity.

Key implications:

- Shifting to Gibbs posteriors allowed us to handle non-zero approximation error
- Provides insights into the trade-off between model complexity and fit

Secondi: Introducing $C(\varepsilon)$

Our oracle inequality doesn't tell us how risk behaves as we get more data.

We propose a new measure of data complexity based on teacher network size.

Secondi: Introducing $C(\varepsilon)$

Our oracle inequality doesn't tell us how risk behaves as we get more data.

We propose a new measure of data complexity based on teacher network size.

Defn. Let r^* is the Bayes error rate. Define

$$C(\varepsilon) = \inf_{\theta^* \in \Theta} m(\theta^*)$$
 s.t. $r(\theta^*) - r^* \leq \varepsilon$

Secondi: Introducing $C(\varepsilon)$

Our oracle inequality doesn't tell us how risk behaves as we get more data.

We propose a new measure of data complexity based on teacher network size.

Defn. Let r^* is the Bayes error rate. Define

$$C(\varepsilon) = \inf_{\theta^* \in \Theta} m(\theta^*)$$
 s.t. $r(\theta^*) - r^* \le \varepsilon$

- ullet Reflects the minimal network size needed for a given approximation error arepsilon
- ullet Monotonically increasing in arepsilon
- Example of a rate-distortion function, up to some approximation.

Secondi: A bound with $C(\varepsilon)$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

Secondi: A bound with $C(\varepsilon)$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

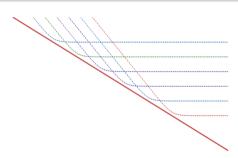
$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$

Secondi: A bound with $C(\varepsilon)$

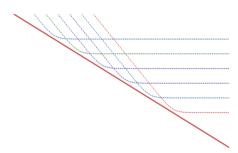
Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$



$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$



$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1} n). \end{split}$$

$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1} n). \end{split}$$

Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then...

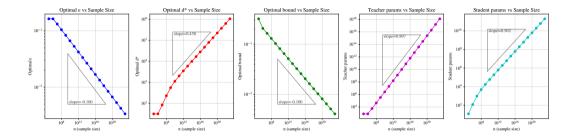
$$r(\hat{\theta}) \leq \inf_{\varepsilon} O\bigg(\varepsilon + \sqrt{\frac{\varepsilon\,C(\varepsilon)\log Q + \log 1/\delta}{n}}\,\bigg)$$

The "complexity" function $C(\varepsilon)$ dictates risk rate:

$$\begin{split} &\text{if } C(\varepsilon) = O(\mathsf{polylog}(\varepsilon^{-1})) & \text{then } r(\hat{\theta}) = O(n^{-1}) \\ &\text{if } C(\varepsilon) = O(\varepsilon^{-p}) & \text{then } r(\hat{\theta}) = O(n^{-1/(p+1)}) \\ &\text{if } C(\varepsilon) = O(\exp(\mathsf{poly}(\varepsilon^{-1}))) & \text{then } r(\hat{\theta}) = O(\log^{-1} n). \end{split}$$

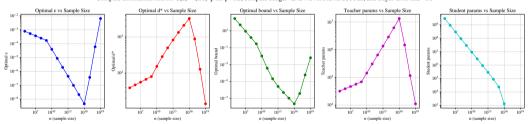
Consequence: If scaling law for n data and size m(n) model says risk decays *slower* than $O(n^{-1/(p+1)})$, then... $C(\varepsilon) \not\in O(\varepsilon^{-p})$.

Scaling laws for $C(\varepsilon) = \varepsilon^{-9}$.

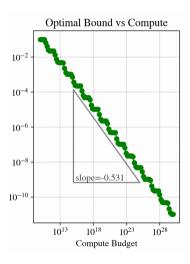


Risk bounds for fixed compute, for $C(\varepsilon) = \varepsilon^{-9}$.





Compute optimal scaling laws for $C(\varepsilon) = \varepsilon^{-9}$.



Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_{n}) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_{n}) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_{n}) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\rho^* \in \Delta(\Theta)} O\bigg(\rho(r) + \sqrt{\frac{KL(\rho||\pi) + \log 1/\delta}{n}}\,\bigg)$$

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_n) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\rho^* \in \Delta(\Theta)} O\bigg(\rho(r) + \sqrt{\frac{KL(\rho||\pi) + \log 1/\delta}{n}}\hspace{1mm}\bigg)$$

Every PAC-Bayes bound (regardless of the predictor) offers a risk bound for the Gibbs posterior-sampled predictor.

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_n) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\rho^* \in \Delta(\Theta)} O\bigg(\rho(r) + \sqrt{\frac{KL(\rho||\pi) + \log 1/\delta}{n}}\,\bigg)$$

- Every PAC-Bayes bound (regardless of the predictor) offers a risk bound for the Gibbs posterior-sampled predictor.
- Offers an explicit connection with coding. Cf. $C(\varepsilon)$ as minimal teacher size.

Donsker-Varadhan offers a more sophisticated lower bound on local entropy.

$$-\lambda^{-1}\log\left(Z_{\lambda}^{\pi}(S)\right) = \inf_{\rho}\rho(\hat{r}_{n}) + \lambda^{-1}KL(\rho||\pi) \tag{5}$$

Theorem

There exists $\lambda > 0$ such that, letting $\hat{\rho}$ be the Gibbs posterior for π , inverse temperature λ , and \hat{r}_n , and letting $\hat{\theta} \mid S \sim \hat{\rho}$, with probability at least $1 - \delta$,

$$r(\hat{\theta}) \leq \inf_{\rho^* \in \Delta(\Theta)} O\bigg(\rho(r) + \sqrt{\frac{KL(\rho||\pi) + \log 1/\delta}{n}}\,\bigg)$$

- Every PAC-Bayes bound (regardless of the predictor) offers a risk bound for the Gibbs posterior-sampled predictor.
- ullet Offers an explicit connection with coding. Cf. C(arepsilon) as minimal teacher size.
- For analysis, we could potentially look to Hessians/flatness (Yang, Mao, Chaudhari 2022). Cf. Dziugaite & Roy 2017.

Key Insights

Our work bridges theory and practice in deep learning generalization.

Main contributions:

- Theoretical explanation for generalization in (mildly) overparameterized models
- Novel measure of data complexity via teacher network size
- Insights into the role of data complexity in learning
- Numerical scaling laws offer us evidence of complexity via excess risk upper bounds

Limitations

Our work has some limitations that offer opportunities for future research.

Current limitations:

- Naive approach to lower bounding probability of interpolation
- Bounds tightest when student size doesn't exceed teacher size

Conclusion

We've introduced a novel perspective on data complexity in deep learning.

Key contributions:

- Novel measure of data complexity based on teacher network size
- Non-vacuous generalization bounds for neural networks
- Connection between theoretical bounds and empirical scaling laws