

Lecture 6: PAC-Bayesian Theory

Lecturer: Liwei Wang

Scribe: Binghui Li, Shiji Xin, Fang Sun, Qizhe Zhang, Mi Yan

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

6.1 Review

VC theory focuses on the uniform convergence for all classifiers in hypothesis space \mathcal{F} . Specifically, with probability $1 - \delta$,

$$P_D \leq P_S + O\left(\sqrt{\frac{d \ln n + \ln 1/\delta}{n}}\right) \quad (6.1)$$

Where P_D is the test error w.r.t. $f \in \mathcal{F}$, P_S is the training error w.r.t. $f \in \mathcal{F}$, d is the VC-dim of \mathcal{F} .

Notice that the result is independent of the distribution D and the learning algorithm A , and only depends on the hypothesis space \mathcal{F} and sample size n .

However, VC Theory fails to account for the great generalization ability of neural networks. For a neural network, the number of parameters is usually far greater than that of training data. In such a case, the generalization bound yielded by VC Theory is scarcely meaningful. Perhaps randomness plays a vital role behind the huge success of neural networks. After all, stochastic gradient descent (SGD), rather than gradient descent (GD), is the common method for optimizing neural networks. [1-3]

6.2 Comparison of Frequentist and Bayesian

There are two points of view of statistical inference: Frequentist and Bayesian. Here is a brief comparison:

	Frequentist	Bayesian
Views of Probability	Law of Large Numbers	Degree of Belief
Parameters Estimation	Maximum Likelihood Estimate	Maximum a posteriori estimation $P(\theta x) = \frac{P(x \theta)P(\theta)}{P(x)}$
Output of Learning	a classifier f	a distribution of classifiers \mathcal{Q}
Prior	Hypothesis space \mathcal{F}	distribution of classifiers \mathcal{P}
Performance	$\text{Err}_D(f)$	$\mathbb{E}_{f \sim \mathcal{Q}}[\text{Err}_D(f)]$
Generalization	VC Theory uniform convergence for all classifiers in a hypothesis space	PAC-Bayesian Theory for all distributions of classifiers
Gap between $\text{Err}_D(f)$ and $\text{Err}_S(f)$	$O\left(\sqrt{\frac{d \ln n + \ln(1/\delta)}{n}}\right)$	$O\left(\sqrt{\frac{D(\mathcal{Q} \mathcal{P}) + \ln(3/\delta)}{n}}\right)$

In the Frequentist view, the gap between test error and training error only depends on hypothesis space \mathcal{F} and sample size, i.e. w/ prob $1 - \delta$, $\forall h \in \mathcal{F}$,

$$\text{Err}_D(h) \leq \text{Err}_S(h) + O\left(\sqrt{\frac{d \ln n + \ln(1/\delta)}{n}}\right)$$

Where $d = \text{VCD}(\mathcal{F})$.

6.3 The main theorem of PAC-Bayesian Theory

In this section, we are going to present the main theorem in PAC-Bayesian Theory.

Theorem 6.1 (PAC-Bayesian) *For any fixed prior distribution of classifiers \mathcal{P} , with probability $1 - \delta$ over the random draw of training dataset S of size n , the following inequality holds uniformly for all distributions (i.e. stochastic classifier) \mathcal{Q} :*

$$\mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)] \leq \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_S(h)] + \sqrt{\frac{D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) + \ln(3/\delta)}{n}} \quad (6.2)$$

Here we denote $\mathbb{P}_{(x,y) \sim D}[y \neq h(x)]$ by $\text{Err}_D(h)$ and $\mathbb{P}_{(x,y) \sim S}[y \neq h(x)]$ by $\text{Err}_S(h)$, $D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) := \mathbb{E}_{h \sim \mathcal{Q}} \left[\ln \frac{q(x)}{p(x)} \right]$ is the KL-divergence between \mathcal{Q} and \mathcal{P} .

Intuitively, since \mathcal{P} is independent of the dataset $S \sim D^n$, we can bound the gap between $\mathbb{E}_{h' \sim \mathcal{P}}[\text{Err}_D(h')]$ and $\mathbb{E}_{h' \sim \mathcal{P}}[\text{Err}_S(h')]$ by Chernoff bound, so we only need to bound the gap between $\mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)]$ and $\mathbb{E}_{h' \sim \mathcal{P}}[\text{Err}_D(h')]$. Before proving the theorem, we list some useful lemmas here.

Lemma 6.2 (Change of Measure) *For all distribution \mathcal{P}, \mathcal{Q} over hypothesis space \mathcal{F} and all function $f : \mathcal{F} \rightarrow \mathbb{R}$, we have*

$$\mathbb{E}_{h \sim \mathcal{Q}}[f(h)] \leq \ln \mathbb{E}_{h' \sim \mathcal{P}} \left[e^{f(h')} \right] + D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) \quad (6.3)$$

Proof: In fact,

$$\begin{aligned} \text{RHS} - \text{LHS} &= \ln \mathbb{E}_{h' \sim \mathcal{P}} \left[e^{f(h')} \right] + D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) - \mathbb{E}_{h \sim \mathcal{Q}}[f(h)] \\ &= \mathbb{E}_{h \sim \mathcal{Q}} \left[\ln \frac{q(h)}{p(h)} \right] - \mathbb{E}_{h \sim \mathcal{Q}}[f(h)] + \ln \mathbb{E}_{h' \sim \mathcal{P}} \left[e^{f(h')} \right] \\ &= \mathbb{E}_{h \sim \mathcal{Q}} \left[\ln \frac{q(h)}{\frac{p(h)e^{f(h)}}{\mathbb{E}_{h' \sim \mathcal{P}} e^{f(h')}}} \right] \\ &= D_{\text{KL}} \left(\mathcal{Q} \parallel \frac{p(h)e^{f(h)}}{\mathbb{E}_{h' \sim \mathcal{P}} e^{f(h')}} \right) \geq 0 \end{aligned}$$

■

Lemma 6.3 *For any $\delta > 0$,*

$$\mathbb{P}_{S \sim D^n} \left[\mathbb{E}_{h \sim \mathcal{P}} \left[e^{n(\text{Err}_D(h) - \text{Err}_S(h))^2} \right] \geq 3/\delta \right] \leq \delta \quad (6.4)$$

Proof: We begin our proof by considering a bound for a fixed $h \sim \mathcal{P}$.

$$\mathbb{E}_{S \sim D^n} \left[e^{n(\text{Err}_D(h) - \text{Err}_S(h))^2} \right] \leq 3$$

For simplicity, let $\Delta := |\text{Err}_D(h) - \text{Err}_S(h)|$. Using Chernoff bound, we have

$$\mathbb{P}_{S \sim D^n} [\Delta \geq \varepsilon] \leq 2 \exp(-2n\varepsilon^2)$$

Then,

$$\begin{aligned} \mathbb{E}_{S \sim D^n} [e^{n\Delta^2}] &= \int_0^\infty \mathbb{P}_{S \sim D^n} [e^{n\Delta^2} \geq t] dt \\ &= \int_1^\infty \mathbb{P}_{S \sim D^n} [e^{n\Delta^2} \geq t] dt + \int_0^1 \mathbb{P}_{S \sim D^n} [e^{n\Delta^2} \geq t] dt \\ &= \int_1^\infty \mathbb{P}_{S \sim D^n} \left[\Delta \geq \sqrt{\frac{\ln t}{n}} \right] dt + 1 \\ &\leq \int_1^\infty 2e^{-2 \ln t} dt + 1 \\ &= 3 \end{aligned}$$

By applying Markov's Inequality, we get

$$\begin{aligned} \mathbb{P}_{S \sim D^n} \left[\mathbb{E}_{h \sim \mathcal{P}} \left[e^{n(\text{Err}_D(h) - \text{Err}_S(h))^2} \right] \geq 3/\delta \right] &= \mathbb{P}_{S \sim D^n} \left[\mathbb{E}_{h \sim \mathcal{P}} [e^{n\Delta^2}] \geq 3/\delta \right] \\ &\leq \frac{\mathbb{E}_{S \sim D^n} \left[\mathbb{E}_{h \sim \mathcal{P}} [e^{n\Delta^2}] \right]}{3/\delta} \\ &= \frac{\mathbb{E}_{h \sim \mathcal{P}} \left[\mathbb{E}_{S \sim D^n} [e^{n\Delta^2}] \right]}{3/\delta} \quad (\text{Fubini Theorem}) \\ &\leq \frac{\mathbb{E}_{h \sim \mathcal{P}} [3]}{3/\delta} = \delta \end{aligned}$$

■

Now we can prove PAC-Bayesian Theorem (6.1) by applying Jensen's Inequality.

With probability $1 - \delta$,

$$\begin{aligned} (\mathbb{E}_{h \sim \mathcal{Q}} [\text{Err}_D(h) - \text{Err}_S(h)])^2 &\leq \frac{1}{n} \mathbb{E}_{h \sim \mathcal{Q}} [n(\text{Err}_D(h) - \text{Err}_S(h))^2] \\ &\leq \frac{1}{n} \left(\ln \mathbb{E}_{h' \sim \mathcal{P}} [e^{n\Delta^2}] + D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) \right) \\ &\leq \frac{1}{n} (\ln(3/\delta) + D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P})) \end{aligned}$$

Thus,

$$\mathbb{E}_{h \sim \mathcal{Q}} [\text{Err}_D(h)] \leq \mathbb{E}_{h \sim \mathcal{Q}} [\text{Err}_S(h)] + \sqrt{\frac{D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) + \ln(3/\delta)}{n}}$$

Lemma 6.4 For any fixed $h \sim \mathcal{P}$, we have

$$\mathbb{E}_{S \sim D^n} \left[e^{nD_{\text{B}}(\text{Err}_S(h) \parallel \text{Err}_D(h))} \right] \leq n + 1 \quad (6.5)$$

Proof: For simplicity, we denote $\text{Err}_D(h)$ by p , $\text{Err}_S(h)$ by \hat{p}_S . Since \mathcal{P} is the prior, h is independent of S . Therefore,

$$\begin{aligned}\mathbb{E}_{S \sim D^n} \left[e^{n D_B(\text{Err}_S(h) \| \text{Err}_D(h))} \right] &= \mathbb{E}_{S \sim D^n} \left[e^{n D_B(\hat{p}_S \| p)} \right] \\ &= \mathbb{E}_{S \sim D^n} \left[\left(\frac{\hat{p}_S}{p} \right)^{n \hat{p}_S} \left(\frac{1 - \hat{p}_S}{1 - p} \right)^{n(1 - \hat{p}_S)} \right]\end{aligned}$$

Note that $S \sim D^n$ implies $\sum_{i=1}^n I[y_i \neq h(x_i)] \sim B(n, p)$, so the equation above can be rewritten as:

$$\begin{aligned}\mathbb{E}_{S \sim D^n} \left[\left(\frac{\hat{p}_S}{p} \right)^{n \hat{p}_S} \left(\frac{1 - \hat{p}_S}{1 - p} \right)^{n(1 - \hat{p}_S)} \right] &= \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \left(\frac{k/n}{p} \right)^{n(k/n)} \left(\frac{1 - k/n}{1 - p} \right)^{n(1 - k/n)} \\ &= \sum_{k=0}^n \binom{n}{k} \left(\frac{k}{n} \right)^k \left(1 - \frac{k}{n} \right)^{n-k} \\ &\leq \sum_{k=0}^n \left(\binom{n}{0} \left(\frac{k}{n} \right)^0 \left(1 - \frac{k}{n} \right)^{n-0} + \binom{n}{1} \left(\frac{k}{n} \right)^1 \left(1 - \frac{k}{n} \right)^{n-1} + \dots \right. \\ &\quad \left. + \binom{n}{k} \left(\frac{k}{n} \right)^k \left(1 - \frac{k}{n} \right)^{n-k} + \dots + \binom{n}{n} \left(\frac{k}{n} \right)^n \left(1 - \frac{k}{n} \right)^{n-n} \right) \\ &= \sum_{k=0}^n \sum_{r=0}^n \binom{n}{r} \left(\frac{k}{n} \right)^r \left(1 - \frac{k}{n} \right)^{n-r} \\ &= \sum_{k=0}^n \left(\frac{k}{n} + \left(1 - \frac{k}{n} \right) \right)^n \\ &= n + 1\end{aligned}$$

■

Using Lemma 6.4, we can get a better result:

Theorem 6.5 *With probability $1 - \delta$,*

$$D_B(\mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_S(h)] \| \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)]) \leq \frac{1}{n} \left(D_B(\mathcal{Q} \| \mathcal{P}) + \ln \left(\frac{n+1}{\delta} \right) \right) \quad (6.6)$$

Proof: For fixed $h' \sim \mathcal{P}$, by applying Markov Inequality and Lemma 6.4, we have

$$\begin{aligned}\mathbb{P}_{S \sim D^n} \left[\mathbb{E}_{h' \sim \mathcal{P}} \left[e^{n D_B(\text{Err}_S(h') \| \text{Err}_D(h'))} \right] \geq \frac{n+1}{\delta} \right] &\leq \frac{\mathbb{E}_{h' \sim \mathcal{P}} \left[\mathbb{E}_{S \sim D^n} \left[e^{n D_B(\text{Err}_S(h') \| \text{Err}_D(h'))} \right] \right]}{\frac{n+1}{\delta}} \\ &\leq \frac{n+1}{\frac{n+1}{\delta}} = \delta\end{aligned}$$

Since $D_B(\cdot\|\cdot)$ is convex, by Lemma 6.2 and Jensen's Inequality, with probability $1 - \delta$,

$$\begin{aligned}
D_B(\mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_S(h)] \| \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)]) &\leq \mathbb{E}_{h \sim \mathcal{Q}}[D_B(\text{Err}_S(h) \| \text{Err}_D(h))] \\
&= \frac{1}{n} \mathbb{E}_{h \sim \mathcal{Q}}[n D_B(\text{Err}_S(h) \| \text{Err}_D(h))] \\
&\leq \frac{1}{n} \left(\ln \left(\mathbb{E}_{h' \sim \mathcal{P}} \left[e^{n D_B(\text{Err}_S(h') \| \text{Err}_D(h'))} \right] \right) + D_{\text{KL}}(\mathcal{Q} \| \mathcal{P}) \right) \\
&\leq \frac{1}{n} \left(\ln \left(\frac{n+1}{\delta} \right) + D_{\text{KL}}(\mathcal{Q} \| \mathcal{P}) \right)
\end{aligned}$$

Thus proving the claim. ■

6.4 PAC-Bayesian Bound for SVM

Now we will apply PAC-Bayesian Theory to linear classifiers. Suppose \mathcal{Q} is a distribution over classifiers, let us define a deterministic voting classifier $g_{\mathcal{Q}}(\mathbf{x})$

$$g_{\mathcal{Q}}(\mathbf{x}) = \text{sgn}(\mathbb{E}_{h \sim \mathcal{Q}} h(\mathbf{x})) \quad (6.7)$$

We have the following proposition:

Proposition 6.6

$$\text{Err}_D(g_{\mathcal{Q}}) \leq 2 \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)] \quad (6.8)$$

Proof: For each point $(\mathbf{x}, y) \sim D$, $g_{\mathcal{Q}}(\mathbf{x}) \neq y$ implies that at least half of h s in \mathcal{Q} satisfy $h(\mathbf{x}) \neq y$, which leads to the conclusion that $\text{Err}_D(g_{\mathcal{Q}}) \leq 2 \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)]$. ■

Now we consider linear classifiers $h(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$, $\mathbf{x} \in \mathbb{R}^d$, $y \in \{\pm 1\}$, $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$.

If we assume $\mathcal{P} = U(S^{d-1})$, the spherical integral of $D_{\text{KL}}(\mathcal{Q} \| \mathcal{P})$ would be complicated to compute. Alternatively, we assume $\mathcal{P} = \mathcal{N}(\mathbf{0}, I_d)$, $\mathcal{Q} = \mathcal{N}(\mu \mathbf{w}, I_d)$, here $\|\mathbf{w}\|_2 = 1$ and μ is a scale factor. According to PAC-Bayesian Theorem (6.1), we have

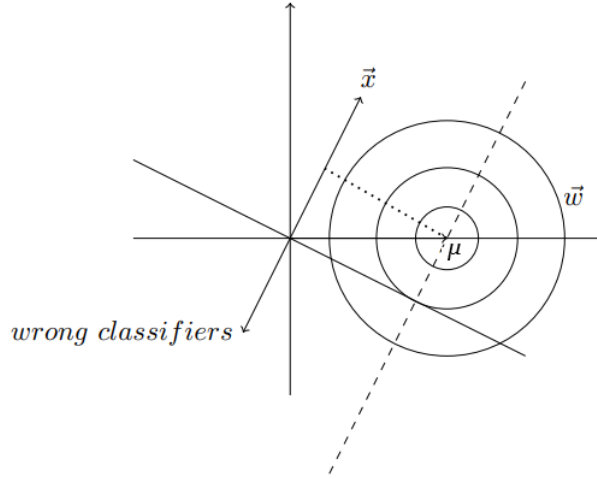
$$\text{Err}_D(g_{\mathcal{Q}}) \leq 2 \left[\text{Err}_S(\mathcal{N}(\mu \mathbf{w}, I_d)) + \sqrt{\frac{D_{\text{KL}}(\mathcal{Q} \| \mathcal{P}) + \ln(3/\delta)}{n}} \right] \quad (6.9)$$

It suffices to determine $D_{\text{KL}}(\mathcal{Q} \| \mathcal{P})$ and $\text{Err}_S(\mathcal{N}(\mu \mathbf{w}, I_d))$ respectively.

6.4.1 Determine $D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P})$

$$\begin{aligned}
 D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) &= \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp \left[-\frac{1}{2} \|\mathbf{x} - \mu \mathbf{w}\|^2 \right] \frac{1}{2} \left(\|\mathbf{x}\|^2 - \|\mathbf{x} - \mu \mathbf{w}\|^2 \right) d\mathbf{x} \\
 &= \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp \left[-\frac{1}{2} \|\mathbf{x} - \mu \mathbf{w}\|^2 \right] \left(\mu \mathbf{w}^\top \mathbf{x} - \frac{1}{2} \mu^2 \right) d\mathbf{x} \\
 &= -\frac{1}{2} \mu^2 + \int_{\mathbb{R}^d} \frac{1}{(\sqrt{2\pi})^d} \exp \left[-\frac{1}{2} \|\mathbf{x} - \mu \mathbf{w}\|^2 \right] (\mu \mathbf{w}^\top (\mathbf{x} - \mu \mathbf{w}) + \mu^2) d\mathbf{x} \\
 &= \frac{1}{2} \mu^2 + \mu \int_{\mathbb{R}^d} (\mathbf{w}^\top \mathbf{x}) \frac{1}{(\sqrt{2\pi})^d} \exp \left[-\frac{1}{2} \|\mathbf{x}\|^2 \right] d\mathbf{x} \\
 &= \frac{1}{2} \mu^2 + \mu \mathbf{w}^\top \mathbb{E}_{\mathbf{x} \sim \mathcal{P}}[\mathbf{x}] \\
 &= \frac{1}{2} \mu^2
 \end{aligned}$$

6.4.2 Determine $\text{Err}_S(\mathcal{N}(\mu \mathbf{w}, I_d))$



For a fixed point (\mathbf{x}, y) and a classifier $\mathbf{h} \sim \mathcal{N}(\mu \mathbf{w}, I_d)$, with the intuition from the figure above, we have

$$\frac{y \mathbf{x}^\top \mathbf{h}}{\|\mathbf{x}\|} \sim \mathcal{N} \left(\frac{y \mathbf{x}^\top (\mu \mathbf{w})}{\|\mathbf{x}\|}, \frac{y^2 \mathbf{x}^\top \mathbf{x}}{\|\mathbf{x}\|^2} \right) = \mathcal{N} \left(\frac{y \mu (\mathbf{w}^\top \mathbf{x})}{\|\mathbf{x}\|}, 1 \right)$$

So

$$\begin{aligned}
 \text{Err}_S(\mathcal{N}(\mu \mathbf{w}, I_d)) &= \mathbb{P}_{\mathbf{h} \sim \mathcal{N}(\mu \mathbf{w}, I_d)} \left[\frac{y (\mathbf{x}^\top \mathbf{h})}{\|\mathbf{x}\|} < 0 \right] \\
 &= \bar{\Phi} \left(\frac{y \mu (\mathbf{w}^\top \mathbf{x})}{\|\mathbf{x}\|} \right)
 \end{aligned}$$

Here, $\Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$, $\bar{\Phi}(t) = 1 - \Phi(t)$.

6.4.3 Putting them together

By PAC-Bayesian Theorem (6.1), for SVM model

$$\begin{aligned} \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_D(h)] &\leq \mathbb{E}_{h \sim \mathcal{Q}}[\text{Err}_S(h)] + \sqrt{\frac{D_{\text{KL}}(\mathcal{Q} \parallel \mathcal{P}) + \ln \frac{3}{\delta}}{n}} \\ &= \mathbb{E}_{S \sim D^n} \left[\bar{\Phi} \left(\frac{y\mu(\mathbf{w}^\top \mathbf{x})}{\|\mathbf{x}\|} \right) \right] + \sqrt{\frac{\frac{\mu^2}{2} + \ln \frac{3}{\delta}}{n}} \end{aligned}$$

Therefore, with probability at least $1 - \delta$ over the random draw of n training data, for all μ and \mathbf{w} , we have

$$\text{Err}_D(g_{\mathcal{Q}}) \leq 2 \left[\mathbb{E}_{S \sim D^n} \left[\bar{\Phi} \left(\frac{y\mu(\mathbf{w}^\top \mathbf{x})}{\|\mathbf{x}\|} \right) \right] + \sqrt{\frac{\frac{\mu^2}{2} + \ln \frac{3}{\delta}}{n}} \right]$$

By optimizing the value of μ , we have

$$\text{Err}_D(g_{\mathcal{Q}}) \leq 2 \inf_{\mu > 0} \left\{ \left[\mathbb{E}_{S \sim D^n} \left[\bar{\Phi} \left(\frac{y\mu(\mathbf{w}^\top \mathbf{x})}{\|\mathbf{x}\|} \right) \right] + \sqrt{\frac{\frac{\mu^2}{2} + \ln \frac{3}{\delta}}{n}} \right] \right\}$$

Since the Gaussian tail probability is monotonically decreasing w.r.t. μ and $\mu^2/2n$ is monotonically increasing w.r.t. μ , we cannot determine the monotonicity of the bound. Intuitively, the Gaussian tail probability will be small if the margin is large, thus the larger the margin is, the tighter the bound will be, which means the generalization ability of SVM will be good when it has a large margin.

References

- [1] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 1225–1234, 2016.
- [2] Samuel Smith, Erich Elsen, and Soham De. On the generalization benefit of noise in stochastic gradient descent. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 9058–9067, 2020.
- [3] Samuel L. Smith and Quoc V. Le. A bayesian perspective on generalization and stochastic gradient descent. In *International Conference on Learning Representations*, 2018.