

A Technical Results

We now give detailed proofs of the theorems in the paper.

A.1 Altitude Training Phenomeon

We begin with a proof of our main generalization bound result, namely Theorem 1. The proof is built on top of the following Berry-Esseen type result.

Lemma 5. *Let Z_1, \dots, Z_d be independent Poisson random variables with means $\lambda_j \in \mathbb{R}_+$, and let*

$$S = \sum_{j=1}^d w_j Z_j, \quad \mu = \mathbb{E}[S], \quad \text{and} \quad \sigma^2 = \text{Var}[S]$$

for some fixed set of weights $\{w_j\}_{j=1}^d$. Then, writing F_S for the distribution function of S and Φ for the standard Gaussian distribution,

$$\sup_{x \in \mathbb{R}} \left| F_S(x) - \Phi\left(\frac{x - \mu}{\sigma}\right) \right| \leq C_{BE} \sqrt{\frac{\max_j \{w_j^2\}}{\sum_{j=1}^d \lambda_j w_j^2}}, \quad (20)$$

where $C_{BE} \leq 4$.

Proof. Our first step is to write S as a sum of bounded *i.i.d.* random variables. Let $N = \sum_{j=1}^d Z_j$. Conditional on N , the Z_j are distributed as a multinomial with parameters $\pi_j = \lambda_j / \lambda$ where $\lambda = \sum_{j=1}^d \lambda_j$. Thus,

$$\mathcal{L}(S | N) \stackrel{d}{=} \mathcal{L}\left(\sum_{k=1}^N W_k | N\right),$$

where $W_k \in \{w_1, \dots, w_d\}$ is a single multinomial draw from the available weights with probability parameters $\mathbb{P}[W_k = w_j] = \pi_j$. This implies that,

$$S \stackrel{d}{=} \sum_{k=1}^N W_k,$$

where N itself is a Poisson random variable with mean λ .

We also know that a Poisson random variable can be written as a limiting mixture of many rare Bernoulli trials:

$$B^{(m)} \Rightarrow N, \quad \text{with} \quad B^{(m)} = \text{Binom}\left(m, \frac{\lambda}{m}\right).$$

The upshot is that

$$S^{(m)} \Rightarrow S, \quad \text{with} \quad S^{(m)} = \sum_{k=1}^m W_k I_k, \quad (21)$$

where the W_k are as before, and the I_k are independent Bernoulli draws with parameter λ/m . Because $S^{(m)}$ converges to S in distribution, it suffices to show that (20) holds for large enough m . The moments of $S^{(m)}$ are correct in finite samples: $\mathbb{E}[S^{(m)}] = \mu$ and $\text{Var}[S^{(m)}] = \sigma^2$ for all m .

The key ingredient in establishing (20) is the Berry-Esseen inequality [see, e.g., 26], which in our case implies that

$$\sup_{x \in \mathbb{R}} \left| F_{S^{(m)}}(x) - \Phi\left(\frac{x - \mu}{\sigma}\right) \right| \leq \frac{\rho_m}{2s_m^3 \sqrt{m}},$$

where

$$s_m^2 = \text{Var}[W_k I_k],$$

$$\rho_m = \mathbb{E}\left[|W_k I_k - \mathbb{E}[W_k I_k]|^3\right],$$

We can show that

$$s_m^2 = \mathbb{E} \left[(W_k I_k)^2 \right] - \mathbb{E} [W_k I_k]^2 = \frac{\lambda}{m} \mathbb{E} [W_k^2] - \left(\frac{\lambda}{m} \mathbb{E} [W_k] \right)^2, \text{ and}$$

$$\rho_m \leq 8 \left(\mathbb{E} [|W_k I_k|^3] + \mathbb{E} [|W_k I_k|^3] \right) = 8 \left(\frac{\lambda}{m} \mathbb{E} [|W_k|^3] + \left(\frac{\lambda}{m} \mathbb{E} [|W_k|] \right)^3 \right).$$

Taking m to ∞ , this implies that

$$\sup_{x \in \mathbb{R}} \left| F_S(x) - \Phi \left(\frac{x - \mu}{\sigma} \right) \right| \leq \frac{4 \mathbb{E} [|W|^3]}{\mathbb{E} [W^2]^{3/2}} \frac{1}{\sqrt{\lambda}}.$$

Thus, to establish (20), it only remains to bound $\mathbb{E} [|W|^3] / \mathbb{E} [W^2]^{3/2}$. Notice that $P_j \stackrel{\text{def}}{=} \pi_j w_j^2 / \mathbb{E} [W^2]$ defines a probability distribution on $\{1, \dots, d\}$, and

$$\frac{\mathbb{E} [|W|^3]}{\mathbb{E} [W^2]} = \mathbb{E}_P [|W| \leq \max_j \{w_j\}].$$

Thus,

$$\frac{\mathbb{E} [|W|^3]}{\mathbb{E} [W^2]^{3/2}} \leq \sqrt{\frac{\max_j \{w_j^2\}}{\sum_{j=1}^d \pi_j w_j^2}}.$$

□

We are now ready to prove our main result.

Proof of Theorem 1. The classifier h is a linear classifier of the form

$$h(x) = \mathbb{I} \{ S > 0 \} \text{ where } S \stackrel{\text{def}}{=} \sum_{j=1}^d w_j x_j,$$

where by assumption $x_j \sim \text{Poisson}(\lambda_j^{(\tau)})$. Our model was fit by dropout, so during training we only get to work with \tilde{x} instead of x , where

$$\tilde{x}_j \sim \text{Binom}(x_j, 1 - \delta), \text{ and so unconditionally}$$

$$\tilde{x}_j \sim \text{Poisson} \left((1 - \delta) \lambda_j^{(\tau)} \right).$$

Without loss of generality, suppose that $c_\tau = 1$, so that we can write the error rate ε_τ during dropout as

$$\varepsilon_\tau = \mathbb{P} \left[\tilde{S} < 0 \mid \tau \right], \text{ where } \tilde{S} = \sum_{j=1}^d w_j \tilde{x}_j. \quad (22)$$

In order to prove our result, we need to translate the information about \tilde{S} into information about S .

The key to the proof is to show that the sums S and \tilde{S} have nearly Gaussian distributions. Let

$$\mu = \sum_{j=1}^d \lambda_j^{(\tau)} w_j \text{ and } \sigma^2 = \sum_{j=1}^d \lambda_j^{(\tau)} w_j^2$$

be the mean and variance of S . After dropout,

$$\mathbb{E} [\tilde{S}] = (1 - \delta) \mu \text{ and } \text{Var} [\tilde{S}] = (1 - \delta) \sigma^2.$$

Writing F_S and $F_{\tilde{S}}$ for the distributions of S and \tilde{S} , we see from Lemma 5 that

$$\begin{aligned} \sup_{x \in \mathbb{R}} \left| F_S(x) - \Phi \left(\frac{x - \mu}{\sigma} \right) \right| &\leq C_{\text{BE}} \sqrt{\Psi_\tau} \text{ and} \\ \sup_{x \in \mathbb{R}} \left| F_{\tilde{S}}(x) - \Phi \left(\frac{x - (1 - \delta)\mu}{\sqrt{1 - \delta}\sigma} \right) \right| &\leq \frac{C_{\text{BE}}}{\sqrt{1 - \delta}} \sqrt{\Psi_\tau}, \end{aligned}$$

where Ψ_τ is as defined in (9). Recall that our objective is to bound $\varepsilon_\tau = F_S(0)$ in terms of $\tilde{\varepsilon}_\tau = F_{\tilde{S}}(0)$. The above result implies that

$$\begin{aligned} \varepsilon_\tau &\leq \Phi \left(-\frac{\mu}{\sigma} \right) + C_{\text{BE}} \sqrt{\Psi_\tau}, \text{ and} \\ \Phi \left(-\sqrt{1 - \delta} \frac{\mu}{\sigma} \right) &\leq \tilde{\varepsilon}_\tau + \frac{C_{\text{BE}}}{\sqrt{1 - \delta}} \sqrt{\Psi_\tau}. \end{aligned}$$

Now, writing $t = \sqrt{1 - \delta} \mu / \sigma$, we can use the Gaussian tail inequalities

$$\frac{\tau}{\tau^2 + 1} < \sqrt{2\pi} e^{\frac{\tau^2}{2}} \Phi(-\tau) < \frac{1}{\tau} \text{ for all } \tau > 0 \quad (23)$$

to check that for all $t \geq 1$,

$$\begin{aligned} \Phi \left(-\frac{t}{\sqrt{1 - \delta}} \right) &\leq \frac{1}{\sqrt{2\pi}} \frac{\sqrt{1 - \delta}}{t} e^{-\frac{t^2}{2(1 - \delta)}} \\ &= \frac{\sqrt{1 - \delta} t^{\frac{\delta}{1 - \delta}}}{\sqrt{2\pi}^{-\frac{\delta}{1 - \delta}}} \left(\frac{1}{\sqrt{2\pi}} \frac{1}{t} e^{-\frac{t^2}{2}} \right)^{\frac{1}{1 - \delta}} \\ &\leq 2^{\frac{1}{1 - \delta}} \frac{\sqrt{1 - \delta} t^{\frac{\delta}{1 - \delta}}}{\sqrt{2\pi}^{-\frac{\delta}{1 - \delta}}} \left(\frac{1}{\sqrt{2\pi}} \frac{t}{t^2 + 1} e^{-\frac{t^2}{2}} \right)^{\frac{1}{1 - \delta}} \\ &\leq \frac{2^{\frac{1}{1 - \delta}} \sqrt{1 - \delta}}{\sqrt{2\pi}^{-\frac{\delta}{1 - \delta}}} t^{\frac{\delta}{1 - \delta}} \Phi(-t)^{\frac{1}{1 - \delta}} \end{aligned}$$

and so noting that in $t\Phi(-t)$ is monotone decreasing in our range of interest and that $t \leq \sqrt{-2 \log \Phi(-t)}$, we conclude that for all $\tilde{\varepsilon}_\tau + C_{\text{BE}}/\sqrt{1 - \delta} \sqrt{\Psi_\tau} \leq \Phi(-1)$,

$$\begin{aligned} \varepsilon_\tau &\leq \frac{2^{\frac{1}{1 - \delta}} \sqrt{1 - \delta}}{\sqrt{4\pi}^{-\frac{\delta}{1 - \delta}}} \left(\sqrt{-\log \left(\tilde{\varepsilon}_\tau + \frac{C_{\text{BE}}}{\sqrt{1 - \delta}} \sqrt{\Psi_\tau} \right)} \right)^{\frac{\delta}{1 - \delta}} \\ &\quad \cdot \left(\tilde{\varepsilon}_\tau + \frac{C_{\text{BE}}}{\sqrt{1 - \delta}} \sqrt{\Psi_\tau} \right)^{\frac{1}{1 - \delta}} + C_{\text{BE}} \sqrt{\Psi_\tau}. \end{aligned} \quad (24)$$

We can also write the above expression in more condensed form:

$$\begin{aligned} &\mathbb{P} \left[\mathbb{I}\{\hat{w} \cdot x^{(i)}\} \neq c_\tau \mid \tau^{(i)} = \tau \right] \\ &= \mathcal{O} \left(\left(\tilde{\varepsilon}_\tau + \sqrt{\frac{\max\{w_j^2\}}{\sum_{j=1}^d \lambda_j^{(\tau)} w_j^2}} \right)^{(1 - \delta)^{\frac{1}{1 - \delta}}} \cdot \max \left\{ 1, \sqrt{-\log(\tilde{\varepsilon}_\tau)^{\frac{\delta}{1 - \delta}}} \right\} \right). \end{aligned} \quad (25)$$

The desired conclusion (9) is equivalent to the above expression, except it uses notation that hides the log factors. \square

Proof of Theorem 2. We can write the dropout error rate as

$$\text{Err}_\delta(\hat{h}_\delta) = \text{Err}_{\min} + \Delta,$$

where Err_{\min} is the minimal possible error from assumption (14) and Δ is the the excess error

$$\Delta = \sum_{\tau=1}^T \mathbb{P}[\tau] \tilde{\varepsilon}_{\tau} \cdot \left| \mathbb{P} \left[y^{(i)} = 1 \mid \tau^{(i)} = \tau \right] - \mathbb{P} \left[y^{(i)} = 0 \mid \tau^{(i)} = \tau \right] \right|.$$

Here, $\mathbb{P}[\tau]$ is the probability of observing a document with topic τ and $\tilde{\varepsilon}_{\tau}$ is as in Theorem 1. The equality follows by noting that, for each topic τ , the excess error rate is given by the rate at which we make sub-optimal guesses, i.e., $\tilde{\varepsilon}_{\tau}$, times the excess probability that we make a classification error given that we made a sub-optimal guess, i.e., $|\mathbb{P}[y^{(i)} = 1 \mid \tau^{(i)} = \tau] - \mathbb{P}[y^{(i)} = 0 \mid \tau^{(i)} = \tau]|$.

Now, thanks to (14), we know that

$$\text{Err}_{\delta}(h_{\delta}^*) = \text{Err}_{\min} + \mathcal{O}\left(\frac{1}{\sqrt{\lambda}}\right),$$

and so the generalization error $\tilde{\eta}$ under the dropout measure satisfies

$$\Delta = \tilde{\eta} + \mathcal{O}\left(\frac{1}{\sqrt{\lambda}}\right).$$

Using (12), we see that

$$\tilde{\varepsilon}_{\tau} \leq \Delta / (2\alpha p_{\min})$$

for each τ , and so

$$\tilde{\varepsilon}_{\tau} = \mathcal{O}\left(\tilde{\eta} + \frac{1}{\sqrt{\lambda}}\right)$$

uniformly in τ . Thus, given the bound (11), we conclude using (25) that

$$\varepsilon_{\tau} = \mathcal{O}\left(\left(\tilde{\eta} + \lambda^{-\frac{1-\delta}{2}}\right)^{\frac{1}{1-\delta}} \max\left\{1, \sqrt{-\log(\tilde{\eta})^{\frac{\delta}{1-\delta}}}\right\}\right)$$

for each topic τ , and so

$$\begin{aligned} \eta &= \text{Err}\left(\hat{h}_{\delta}\right) - \text{Err}\left(h_{\delta}^*\right) \\ &= \mathcal{O}\left(\left(\tilde{\eta} + \lambda^{-\frac{1-\delta}{2}}\right)^{\frac{1}{1-\delta}} \max\left\{1, \sqrt{-\log(\tilde{\eta})^{\frac{\delta}{1-\delta}}}\right\}\right), \end{aligned} \quad (26)$$

which directly implies (16). Note η will in general be larger than the ε_{τ} , because guessing the optimal label c_{τ} is not guaranteed to lead to a correct classification decision (unless each topic is pure, i.e., only represents one class). Here, subtracting the optimal error $\text{Err}(h_{\delta}^*)$ allows us to compensate for this effect. \square

Proof of Corollary 3. Here, we prove the more precise bound

$$\text{Err}\left(\hat{h}_{\delta}\right) - \text{Err}\left(h_{\delta}^*\right) = \mathcal{O}_P\left(\sqrt{\left(\frac{d}{n} + \frac{1}{\lambda^{(1-\delta)}}\right) \max\left\{1, \log\left(\frac{n}{d}\right)\right\}^{1+\delta} \frac{1}{1-\delta}}\right). \quad (27)$$

To do this, we only need to show that

$$\text{Err}_{\delta}\left(\hat{h}_{\delta}\right) - \text{Err}_{\delta}\left(h_{\delta}^*\right) = \mathcal{O}_P\left(\sqrt{\frac{d}{n} \max\left\{1, \log\left(\frac{n}{d}\right)\right\}}\right), \quad (28)$$

i.e., that dropout generalizes at the usual rate with respect to the dropout measure. Then, by applying (26) from the proof of Theorem 2, we immediately conclude that \hat{h}_{δ} converges at the rate given in (17) under the data-generating measure.

Let $\widehat{\text{Err}}_{\delta}(h)$ be the average training loss for a classifier h . The empirical loss is unbiased, i.e.,

$$\mathbb{E}\left[\widehat{\text{Err}}_{\delta}(h)\right] = \text{Err}_{\delta}(h).$$

Given this unbiasedness condition, standard methods for establishing rates as in (28) [e.g., 27] only require that the loss due to any single training example $(x^{(i)}, y^{(i)})$ is bounded, and that the training examples are independent; these conditions are needed for an application of Hoeffding's inequality. Both of these conditions hold here. \square

A.2 Distinct Topics Assumption

Proposition 6. *Let the generative model from Section 2 hold, and define*

$$\pi^{(\tau)} = \lambda^{(\tau)} / \left\| \lambda^{(\tau)} \right\|_1$$

for the topic-wise word probability vectors and

$$\Pi = (\pi^{(1)}, \dots, \pi^{(T)}) \in \mathbb{R}^{d \times T}$$

for the induced matrix. Suppose that Π has rank T , and that the minimum singular value of Π (in absolute value) is bounded below by

$$|\sigma_{\min}(\Pi)| \geq \sqrt{\frac{T}{(1-\delta)\lambda}} \left(1 + \sqrt{\log_+ \frac{\lambda}{2\pi}} \right), \quad (29)$$

where \log_+ is the positive part of \log . Then (14) holds.

Proof. Our proof has two parts. We begin by showing that, given (29), there is a vector w with $\|w\|_2 \leq 1$ such that

$$\mathbb{I} \left\{ w \cdot \pi^{(\tau)} > 0 \right\} = c_\tau, \text{ and } \left| w \cdot \pi^{(\tau)} \right| \geq -\frac{1}{\sqrt{(1-\delta)\lambda}} \Phi^{-1} \left(\frac{1}{\sqrt{\lambda}} \right) \quad (30)$$

for all topics τ ; in other words, the topic centers can be separated with a large margin. After that, we show that (30) implies (14).

We can re-write the condition (30) as

$$\min \left\{ \|w\|_2 : c_\tau w \cdot \pi^{(\tau)} \geq 1 \text{ for all } \tau \right\} \leq \left(-\frac{1}{\sqrt{(1-\delta)\lambda}} \Phi^{-1} \left(\frac{1}{\sqrt{\lambda}} \right) \right)^{-1},$$

or equivalently that

$$\min \left\{ \|w\|_2 : S \Pi^\top w \geq 1 \right\} \leq \left(-\frac{1}{\sqrt{(1-\delta)\lambda}} \Phi^{-1} \left(\frac{1}{\sqrt{\lambda}} \right) \right)^{-1}$$

where $S = \text{diag}(c_\tau)$ is a diagonal matrix of class signs. Now, assuming that $\text{rank}(\Pi) \geq T$, we can verify that

$$\begin{aligned} \min \left\{ \|w\|_2 : S \Pi^\top w \geq 1 \right\} &= \min \left\{ \sqrt{z^\top (\Pi^\top S^2 \Pi)^{-1} z} : z \geq 1 \right\} \\ &\leq \sqrt{1^\top (\Pi^\top \Pi)^{-1} 1} \\ &\leq |\sigma_{\min}(\Pi)|^{-1} \sqrt{T} \\ &\leq \left(\frac{1}{\sqrt{(1-\delta)\lambda}} \left(1 + \sqrt{\log_+ \frac{\lambda}{2\pi}} \right) \right)^{-1}, \end{aligned}$$

where the last line followed by hypothesis. Now, by (23)

$$\Phi \left(-\left(1 + \sqrt{\log_+ \frac{\lambda}{2\pi}} \right) \right) \leq \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} \log \frac{\lambda}{2\pi} \right) = \frac{1}{\sqrt{\lambda}}.$$

Because Φ^{-1} is monotone increasing, this implies that

$$\left(1 + \sqrt{\log_+ \frac{\lambda}{2\pi}} \right)^{-1} \leq \left(-\Phi^{-1} \left(\frac{1}{\sqrt{\lambda}} \right) \right)^{-1},$$

and so (30) holds.

Now, taking (30) as given, it suffices to check that the sub-optimal prediction rate is $\mathcal{O}\left(1/\sqrt{\lambda}\right)$ uniformly for each τ . Focusing now on a single topic τ , suppose without loss of generality that $c_\tau = 1$. We thus need to show that

$$\mathbb{P}[w \cdot \tilde{x} \leq 0] = \mathcal{O}\left(\frac{1}{\sqrt{\lambda}}\right),$$

where \tilde{x} is a feature vector thinned by dropout. By Lemma 5 together with (11), we know that

$$\mathbb{P}[w \cdot \tilde{x} \leq 0] \leq \Phi\left(-\frac{\mathbb{E}[w \cdot \tilde{x}]}{\sqrt{\text{Var}[w \cdot \tilde{x}]}}\right) + \mathcal{O}\left(\frac{1}{\sqrt{\lambda}}\right).$$

By hypothesis,

$$\mathbb{E}[w \cdot \tilde{x}] \geq -\sqrt{(1-\delta)\lambda^{(\tau)}}\Phi^{-1}\left(\frac{1}{\sqrt{\lambda}}\right),$$

and we can check that

$$\text{Var}[w \cdot \tilde{x}] = (1-\delta)\sum_{j=1}^d w_j^2 \lambda_j^{(\tau)} \leq (1-\delta)\lambda^{(\tau)}$$

because $\|w\|_2 \leq 1$. Thus,

$$\Phi\left(-\frac{\mathbb{E}[w \cdot \tilde{x}]}{\sqrt{\text{Var}[w \cdot \tilde{x}]}}\right) \leq \Phi\left(\Phi^{-1}\left(\frac{1}{\sqrt{\lambda}}\right)\right) = \frac{1}{\sqrt{\lambda}},$$

and (14) holds. \square

A.3 Dropout Preserves the Bayes Decision Boundary

Proof of Proposition 4. Another way to view our topic model is as follows. For each topic τ , define a distribution over words $\pi^{(\tau)} \in \Delta^{d-1}$: $\pi^{(\tau)} \stackrel{\text{def}}{=} \lambda^{(\tau)} / \|\lambda^{(\tau)}\|_1$. The generative model is equivalent to first drawing the length of the document and then drawing the words from a multinomial:

$$L_i \sim \text{Poisson}\left(\|\lambda^{(\tau)}\|_1\right), \text{ and } x^{(i)} \mid \tau^{(i)}, L_i \sim \text{Multinom}\left(\pi^{(\tau^{(i)})}, L_i\right). \quad (31)$$

Now, write the multinomial probability mass function (31) as

$$\mathbb{P}_m[x; \pi, L] = \frac{L!}{x_1! \cdots x_d!} \pi_1^{x_1} \cdots \pi_d^{x_d}$$

For each label c , define Π_c to be the distribution over the probability vectors induced by the distribution over topics. Note that we could have an infinite number of topics. By Bayes rule,

$$\mathbb{P}[x = v \mid y = c] = \mathbb{P}\left[L = \sum_{j=1}^d v_j\right] \cdot \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_c(\pi), \text{ and}$$

$$\mathbb{P}[y = c \mid x = v] = \frac{\mathbb{P}[c] \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_c(\pi)}{\sum_{c'} \mathbb{P}[c'] \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_{c'}(\pi)}.$$

The key part is that the distribution of L doesn't depend on c , so that when we condition on $x = v$, it cancels. As for the joint distribution of (\tilde{x}, y) , note that, given π and $\tilde{L} = \sum_{j=1}^d \tilde{x}_j$, \tilde{x} is conditionally $\text{Multinom}(\pi, \tilde{L})$. So then

$$\mathbb{P}[\tilde{x} = v \mid y = c] = \mathbb{P}\left[\tilde{L} = \sum_{j=1}^d v_j\right] \cdot \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_c(\pi), \text{ and}$$

$$\mathbb{P}[y = c \mid \tilde{x} = v] = \frac{\mathbb{P}[c] \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_c(\pi)}{\sum_{c'} \mathbb{P}[c'] \int \mathbb{P}_m\left[v; \pi, \sum_{j=1}^d v_j\right] d\Pi_{c'}(\pi)}.$$

In both cases, L and \tilde{L} don't depend on the topic, and when we condition on x and \tilde{x} , we get the same distribution over y . \square