

A Proofs for Section 2

Proof of Proposition 2.5. Let $\iota \stackrel{\text{def}}{=} \pi_1 \times \cdots \times \pi_k$. Since ι is injective by assumption, we have

$$\exp(A(\beta; z)) = \sum_{y \in \mathcal{Y}} \exp(\beta^\top \psi(z, y)) \quad (16)$$

$$= \sum_{p \in \iota(\mathcal{Y})} \exp(\beta^\top \psi(z, \iota^{-1}(p))) \quad (17)$$

$$= \sum_{p \in \iota(\mathcal{Y})} \exp\left(-\sum_{j=1}^k \beta_j \mathbb{I}[\pi_j(f(z)) \neq p_j]\right) \quad (18)$$

$$\leq \sum_{p \in \prod_j \mathcal{Y}_j} \exp\left(-\sum_{j=1}^k \beta_j \mathbb{I}[\pi_j(f(z)) \neq p_j]\right) \quad (19)$$

$$= \prod_{j=1}^k \sum_{p_j \in \mathcal{Y}_j} \exp(-\beta_j \mathbb{I}[\pi_j(f(z)) \neq p_j]) \quad (20)$$

$$= \prod_{j=1}^k (1 + (|\mathcal{Y}_j| - 1) \exp(-\beta_j)), \quad (21)$$

as was to be shown. \square

B Proofs for Section 3

B.1 Effect on loss

Proof of Proposition 3.1. Note that we have

$$L^* = \mathbb{E}_{p^*}[-\log p_{\theta^*}(y | x)] \quad (22)$$

$$= \mathbb{E}_{p^*}[-\log \mathbb{E}_{z \sim p_{\theta^*}}[\mathbb{S}(z, y)]] \quad (23)$$

$$\stackrel{(a)}{\geq} \mathbb{E}_{p^*}[-\log \mathbb{E}_{z \sim p_{\theta^*}}[\exp(\beta^\top \psi(z, y))]] \quad (24)$$

$$= \mathbb{E}_{p^*}[-\log \mathbb{E}_{z \sim p_{\theta^*}}[p_\beta(y | z) \exp(A(\beta))]] \quad (25)$$

$$= \mathbb{E}_{p^*}[-\log p_{\theta^*, \beta}(y | x) - A(\beta)] \quad (26)$$

$$= L(\theta^*, \beta) - A(\beta) \quad (27)$$

$$\stackrel{(b)}{\geq} L(\theta_\beta^*, \beta) - A(\beta). \quad (28)$$

Here (a) follows because $\mathbb{S}(z, y) \leq \exp(\beta^\top \psi(z, y))$, since the latter is non-negative and is 1 when $\mathbb{S}(z, y) = 1$; (b) follows because θ_β^* is the minimizer of $L(\cdot, \beta)$. Continuing:

$$L(\theta_\beta^*, \beta) - A(\beta) = \mathbb{E}_{p^*}[-\log \mathbb{E}_{z \sim p_{\theta_\beta^*}}[\exp(\beta^\top \psi(z, y))]] \quad (29)$$

$$\stackrel{(c)}{\geq} \mathbb{E}_{p^*}[-\log \mathbb{E}_{z \sim p_{\theta_\beta^*}}[\exp(-\beta_{\min}(1 - \mathbb{S}(z, y)))] \quad (30)$$

$$= \mathbb{E}_{p^*}[-\log(p_{\theta_\beta^*}(y | x) + (1 - p_{\theta_\beta^*}(y | x)) \exp(-\beta_{\min}))] \quad (31)$$

$$= \mathbb{E}_{p^*}[-\log(1 - (1 - p_{\theta_\beta^*}(y | x))(1 - \exp(-\beta_{\min})))] \quad (32)$$

$$\stackrel{(d)}{\geq} \mathbb{E}_{p^*}[(1 - p_{\theta_\beta^*}(y | x))(1 - \exp(-\beta_{\min}))]. \quad (33)$$

Again, (c) follows because $\beta^\top \psi(z, y) \leq -\beta_{\min}(1 - \mathbb{S}(z, y))$, and (d) follows because $-\log(1 - x) \geq x$ for $x \leq 1$. Putting these together, we have $L^* \geq (1 - \exp(-\beta_{\min})) \mathbb{E}_{p^*}[1 - p_{\theta_\beta^*}(y | x)]$, which yields the desired result. \square

Proof of Lemma 3.2. We will show a stronger result: any model and relaxation can be slightly modified to cause $\mathbb{E}_{p^*}[p_{\theta_\beta}(y | x)]$ to be zero, in a way that is demonstrated below (though the modified model will no longer be an exponential family).

Given any $\mathbb{S}_{1:k}$, construct a new point z_0 such that $\mathbb{S}_{1:k}(z_0, y) = 1$ for all y , and add a new constraint $\mathbb{S}_0(z, y) = [z \neq z_0]$. Then $\mathbb{S}(z_0, y) = 0$ for all y , so we never want to place mass on z_0 under the unrelaxed supervision. In addition, extend the model family to allow the single additional distribution $p'(z | x) = \mathbb{I}[z = z_0]$.

Now, suppose $\beta_{1:k} = \infty$ and $\beta_0 = \beta_{\min}$. Then, for any θ , we have $L(\theta, \beta) = A(\beta) + L(\theta, \infty)$, since p_θ places no mass on z_0 ; therefore, $L(\theta, \beta) \geq A(\beta) + L^*$ for all θ . On the other hand, we have $L(p', \beta) = A(\beta) + \beta_{\min}$. If $\beta_{\min} < L^*$, we will thus use p' and shift all of the mass to z_0 , thereby placing zero mass on the correct answer. \square

Note that the proof required constructing a ‘‘bad’’ z_0 that satisfied almost all the constraints for many values of y at once. It seems straightforward to avoid this in practice, and so it would be interesting to find assumptions under which we obtain a better relative loss bound than Proposition 3.1.

B.2 Amount of data needed to learn

For the next few derivations we will make extensive use of the relation $\log p_{\theta, \beta}(y | x) = A(\theta, \beta; x, y) - A(\theta; x)$, where $A(\theta, \beta; x, y) \stackrel{\text{def}}{=} \log(\sum_z \exp(\theta^\top \phi(x, z) + \beta^\top \psi(z, y)))$. Note that the preceding definition is consistent with (13) since we assume throughout Section 3 that $\mathbb{T} \equiv 1$. We will also use the following properties of log-partition functions:

$$\nabla_\theta A(\theta, \beta; x, y) = \mathbb{E}_{z \sim p_{\theta, \beta}(\cdot | x, y)}[\phi(x, z)] \quad (34)$$

$$= \frac{\mathbb{E}_{z \sim p_\theta(\cdot | x)}[\phi(x, z) \exp(\beta^\top \psi(z, y))]}{\mathbb{E}_{z \sim p_\theta(\cdot | x)}[\exp(\beta^\top \psi(z, y))]}, \quad (35)$$

$$\nabla_\theta^2 A(\theta, \beta; x, y) = -(\nabla_\theta A)(\nabla_\theta A)^\top + \mathbb{E}_{z \sim p_{\theta, \beta}(\cdot | x, y)}[\phi(x, z) \otimes \phi(x, z)] \quad (36)$$

$$= -(\nabla_\theta A)(\nabla_\theta A)^\top + \frac{\mathbb{E}_{z \sim p_\theta(\cdot | x)}[(\phi(x, z) \otimes \phi(x, z)) \exp(\beta^\top \psi(z, y))]}{\mathbb{E}_{z \sim p_\theta(\cdot | x)}[\exp(\beta^\top \psi(z, y))]} \quad (37)$$

Here we use $\nabla_\theta A$ as short-hand for $\nabla_\theta A(\theta, \beta; x, y)$. These $\nabla_\theta A$ terms will always cancel out in the sequel, so they can be safely ignored. (The cancellation occurs because we always end up subtracting two log-normalization constants, whose gradients must be equal by first-order optimality conditions.) Analogous properties to those above hold for $A(\theta; x)$:

$$\nabla_\theta A(\theta; x) = \mathbb{E}_{z \sim p_\theta(\cdot | x)}[\phi(x, z)], \quad (38)$$

$$\nabla_\theta^2 A(\theta; x) = -(\nabla_\theta A)(\nabla_\theta A)^\top + \mathbb{E}_{z \sim p_\theta(\cdot | x)}[\phi(x, z) \otimes \phi(x, z)]. \quad (39)$$

In this case, $\nabla_\theta A$ is short-hand for $\nabla_\theta A(\theta; x)$.

Proof of (8). We have

$$\mathcal{I}_\infty = \nabla_\theta^2 [-\log p_{\theta^*, \infty}(y | x)] \quad (40)$$

$$= \nabla_\theta^2 [A(\theta^*; x) - A(\theta^*, \infty; x, y)] \quad (41)$$

$$= \mathbb{E}_{\theta^*}[\phi(x, z) \otimes \phi(x, z)] - \frac{\mathbb{E}_{\theta^*}[(\phi(x, z) \otimes \phi(x, z))\mathbb{S}(z, y)]}{\mathbb{E}_{\theta^*}[\mathbb{S}(z, y)]} \quad (42)$$

$$= \mathbb{E}_{\theta^*}[\phi(x, z) \otimes \phi(x, z)] - \mathbb{E}_{\theta^*}[\phi(x, z) \otimes \phi(x, z) | \mathbb{S}] \quad (43)$$

$$= (\mathbb{P}[\neg \mathbb{S}]\mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg \mathbb{S}] + \mathbb{P}[\mathbb{S}]\mathbb{E}_{\theta^*}[\phi \otimes \phi | \mathbb{S}]) - \mathbb{E}_{\theta^*}[\phi \otimes \phi | \mathbb{S}] \quad (44)$$

$$= \mathbb{P}_{\theta^*}[\neg \mathbb{S}] (\mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg \mathbb{S}] - \mathbb{E}_{\theta^*}[\phi \otimes \phi | \mathbb{S}]). \quad (45)$$

The result follows by taking expectations. \square

Proof of (9). We have

$$\mathcal{I}_\beta = \nabla_\theta^2[-\log p_{\theta^*,\beta}(y | x)] \quad (46)$$

$$= \nabla_\theta^2 [A(\theta^*; x) - A(\theta^*, \beta; x, y)] \quad (47)$$

$$= \mathbb{E}_{\theta^*}[\phi(x, z) \otimes \phi(x, z)] - \frac{\mathbb{E}_{\theta^*}[(\phi(x, z) \otimes \phi(x, z)) \exp(\beta^\top \psi)]}{\mathbb{E}_{\theta^*}[\exp(\beta^\top \psi)]} \quad (48)$$

$$= \frac{\mathbb{E}_{\theta^*}[\phi \otimes \phi] \mathbb{E}_{\theta^*}[\exp(\beta^\top \psi)] - \mathbb{E}_{\theta^*}[(\phi \otimes \phi) \exp(\beta^\top \psi)]}{\mathbb{E}_{\theta^*}[\exp(\beta^\top \psi)]} \quad (49)$$

$$= -\frac{\text{Cov}_{\theta^*}[\phi \otimes \phi, \exp(\beta^\top \psi)]}{\mathbb{E}_{\theta^*}[\exp(\beta^\top \psi)]} \quad (50)$$

$$\stackrel{(a)}{=} -\frac{\text{Cov}_{\theta^*}[\phi \otimes \phi, 1 + \beta^\top \psi + \mathcal{O}(\beta^2)]}{\mathbb{E}_{\theta^*}[1 + \mathcal{O}(\beta)]} \quad (51)$$

$$\stackrel{(b)}{=} -\text{Cov}_{\theta^*}[\phi \otimes \phi, \beta^\top \psi] + \mathcal{O}(\beta^2), \quad (52)$$

where in (a) we used $\exp(\beta^\top \psi) = 1 + \beta^\top \psi + \mathcal{O}(\beta^2)$ and in (b) we used $\text{Cov}[\cdot, 1] = 0$. The result again follows by taking expectations. \square

Note: Assuming that $\|\psi\|_1$ is small for most z (as measured by p_{θ^*}), the $\mathcal{O}(\beta^2)$ term is small as long as $\|\beta\|_\infty \ll 1$. This assumption on ψ holds when $\mathbb{P}_{\theta^*}[\mathbb{S}] \approx 1$ (so that $\psi = 0$ most of the time).

Proof of (10). Recall that we are assuming $\beta_j = \beta_{\min}$ for all j , and that the $\neg\mathbb{S}_j$ are all disjoint. In this case, $-\beta^\top \psi$ is equal to β_{\min} if a constraint is violated and 0 if no constraints are violated. We then have

$$\text{Cov}_{\theta^*}[\phi \otimes \phi, -\beta^\top \psi] \quad (53)$$

$$= \beta_{\min} \text{Cov}_{\theta^*}[\phi \otimes \phi, \mathbb{I}[\neg\mathbb{S}]] \quad (54)$$

$$= \beta_{\min} \mathbb{P}[\neg\mathbb{S}] \left(\mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg\mathbb{S}] - \mathbb{E}_{\theta^*}[\phi \otimes \phi] \right) \quad (55)$$

$$= \beta_{\min} \mathbb{P}[\neg\mathbb{S}] \left(\mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg\mathbb{S}] - \mathbb{P}[\neg\mathbb{S}] \mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg\mathbb{S}] - \mathbb{P}[\mathbb{S}] \mathbb{E}_{\theta^*}[\phi \otimes \phi | \mathbb{S}] \right) \quad (56)$$

$$= \beta_{\min} \mathbb{P}[\mathbb{S}] \mathbb{P}[\neg\mathbb{S}] \left(\mathbb{E}_{\theta^*}[\phi \otimes \phi | \neg\mathbb{S}] - \mathbb{E}_{\theta^*}[\phi \otimes \phi | \mathbb{S}] \right), \quad (57)$$

as claimed. \square

B.3 Optimizing β

Proof of Proposition 3.3. We can re-express $\mathbb{E}_{x,y \sim p^*}[-\log p(y | x)]$ as $\text{KL}(p^* \| p) + H(p^*)$. Hence, in particular, $L(\theta, \beta) = \text{KL}(p^* \| p_{\theta,\beta}) + H(p^*) \geq H(p^*)^2$, with equality if and only if $p_{\theta,\beta} = p^*$. On the other hand, $p_{\theta^*,\infty} = p_{\theta^*} = p^*$ by assumption, so equality is attainable, and (θ^*, ∞) is a global optimum of L .

Note that the normalization constant $A(\beta)$ is important here, since if $p_{\theta,\beta}$ did not (sub-)normalize then the KL divergence would not necessarily be non-negative. \square

²Here we use the fact that $\text{KL}(p \| q) \stackrel{\text{def}}{=} \mathbb{E}_p[\log p - \log q]$ is non-negative as long as p normalizes and q sub-normalizes, which is true for $q = p_{\theta,\beta}$ by Proposition 2.5).

C Proofs for Section 4

Proof of (14). The acceptance rate is simply the expectation, over all $z \mid x$, of the acceptance probability for that particular z . This can clearly be written as

$$\sum_z p_{\theta, \mathbb{T}}(z \mid x) \exp(\beta^\top \psi(z, y)) \quad (58)$$

$$= \sum_z \mathbb{T}(z, y) \exp(\theta^\top \phi(x, z) - A_{\mathbb{T}}(\theta; x, y)) \exp(\beta^\top \psi(z, y)) \quad (59)$$

$$= \exp(-A_{\mathbb{T}}(\theta; x, y)) \sum_z \mathbb{T}(z, y) \exp(\theta^\top \phi(x, z) + \beta^\top \psi(z, y)) \quad (60)$$

$$= \exp(A(\theta, \beta; x, y) - A_{\mathbb{T}}(\theta; x, y)). \quad (61)$$

Since (14) is the multiplicative inverse of (61), the result follows. \square

Proof of (15). By convexity of $A(\theta, \beta; x, y)$, we have

$$A(\theta, \beta; x, y) \quad (62)$$

$$\geq A(\tilde{\theta}, \tilde{\beta}; x, y) + (\theta - \tilde{\theta})^\top \nabla_\theta A(\tilde{\theta}, \tilde{\beta}; x, y) + (\beta - \tilde{\beta})^\top \nabla_\beta A(\tilde{\theta}, \tilde{\beta}; x, y) \quad (63)$$

$$= A(\tilde{\theta}, \tilde{\beta}; x, y) + (\theta - \tilde{\theta})^\top \mathbb{E}_{p_{\tilde{\theta}, \tilde{\beta}}(\cdot \mid x, y)}[\phi(x, z)] + (\beta - \tilde{\beta})^\top \mathbb{E}_{p_{\tilde{\theta}, \tilde{\beta}}(\cdot \mid x, y)}[\psi(z, y)] \quad (64)$$

$$= A(\tilde{\theta}, \tilde{\beta}; x, y) + (\theta - \tilde{\theta})^\top \tilde{\phi} + (\beta - \tilde{\beta})^\top \tilde{\psi}, \quad (65)$$

as was to be shown. \square