
Supplementary Material

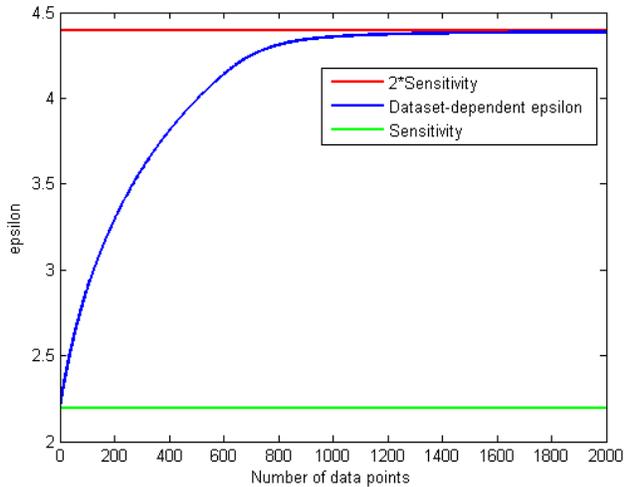


Figure 1: An adversary greedily selects data points to add to a dataset to increase the dataset-specific privacy cost ϵ of posterior sampling via the exponential mechanism (OPS).

A ADVERSARIAL DATA EXPERIMENT

In this appendix we describe an additional simulation experiment which supplements the analysis performed in the main manuscript. Wang et al. (2015)’s analysis finds that the privacy cost of posterior sampling does not directly improve with the number of data points N , unless the analyst deliberately modifies the posterior by changing the temperature before sampling. In Figure 1 we report an experiment showing that this result is not just a limitation of the analysis: there do exist cases where the dataset-specific privacy cost of posterior sampling can approach the exponential mechanism worst case of $\epsilon = 2\Delta \log Pr(\theta, \mathbf{X})$ as the number of observations N increases.

In the experiment, we consider a beta distribution posterior, symmetrically truncated at $a_0 = 0.1$, with Bernoulli observations. We simulate an adversary who greedily selects data points to add to a dataset to increase the dataset-

specific privacy cost ϵ of posterior sampling. The dataset-specific “local” privacy parameter ϵ is computed via a grid search over the Bernoulli success parameter p and Bernoulli outcomes x, x' , for the case where the adversary adds a success, or a failure, and the adversary selects the success/failure outcome with the highest local ϵ . The adversary is able to make the dataset-specific ϵ approach the worst case by manipulating the partition function of the posterior. The exponential mechanism’s worst case for posterior sampling, $\epsilon = 2\Delta \log Pr(\theta, \mathbf{X})$, corresponds to a sum of two cost terms. We must pay a cost of $\Delta \log Pr(\theta, \mathbf{X})$ from to the difference of log-likelihood terms, as we can always draw the worst-case θ (e.g., when p is on the truncation boundary), plus another $\Delta \log Pr(\theta, \mathbf{X})$ in the worst case due to the difference of log partition-functions terms, which the adversary can alter up to the worst case, as they do in Figure 1. This is described formally in the supplementary of (Wang et al., 2015).

B PROOFS OF THEORETICAL RESULTS

Here we provide proofs for the results presented in Section 3.3.

B.1 PROOF OF LAPLACE MECHANISM ASYMPTOTIC KL-DIVERGENCE

Our results hold specifically over the class of exponential families. A family of distributions parameterized by θ which has the form

$$Pr(\mathbf{x}|\theta) = h(\mathbf{x}) \exp\left(\theta^T S(\mathbf{x}) - A(\theta)\right) \quad (1)$$

is said to be an exponential family. Breaking down this structure into its parts, θ is a vector known as the natural parameters for the distribution and lies in some space Θ . $S(\mathbf{x})$ represents a vector of sufficient statistics that fully capture the information needed to determine how likely \mathbf{x} is under this distribution. $A(\theta)$ represents the log-normalizer,

a term used to make this a probability distribution sum to one over all possibilities of \mathbf{x} . $h(\mathbf{x})$ is a base measure for this family, independent of which distribution in the family is used.

As we are interested in learning θ , we are considering algorithms that generate a posterior distribution for θ . The exponential families always have a conjugate prior family which is itself an exponential family. When speaking of these prior and posterior distributions, θ becomes the random variable and we introduce a new vector of natural parameters η in a space M to parameterize these distributions. To ease notation, we will express this conjugate prior exponential family as $Pr(\theta|\eta) = f(\theta) \exp(\eta^\top T(\theta) - B(\eta))$, which is simply a relabelling of the exponential family structure. The posterior from this conjugate prior is often written in an equivalent form

$$Pr(\theta|\mathbf{X}, \chi, \alpha) \propto g(\theta)^{N+\alpha} \exp\left(\theta^\top \left(\sum_{i=1}^N S(\mathbf{x}^{(i)}) + \alpha\chi\right)\right),$$

where the vector χ and the scalar α together specify the vector η of natural parameters for this distribution. From the interaction of χ , α , and \mathbf{X} on the posterior, one can see that this prior acts like α observations with average sufficient statistics χ have already been observed. This parameterization with χ and α has many nice intuitive properties, but our proofs center around the natural parameter vector η for this prior.

These two forms for the posterior can be reconciled by letting $\eta = (\alpha\chi + \sum_{i=1}^N S(\mathbf{x}^{(i)}), N + \alpha)$ and $T(\theta) = (\theta, -A(\theta))$. This definition for the natural parameters η and sufficient statistics $T(\theta)$ fully specify the exponential family the posterior resides in, with $B(\eta)$ defined as the appropriate log-normalizer for this distribution (and $f(\theta) = 1$ is merely a constant). We note that the space of $T(\Theta)$ is not the full space \mathbb{R}^{d+1} , as the last component of $T(\theta)$ is a function of the previous components. Plugging in these expressions for η and $T(\theta)$ we get the following form for the conjugate prior:

$$\begin{aligned} Pr(\theta|\mathbf{X}, \chi, \alpha) &= \exp\left(\theta^\top \left(\alpha\chi + \sum_{i=1}^N S(\mathbf{x}^{(i)})\right)\right) \\ &\quad - (N + \alpha)A(\theta) \\ &\quad - B(\eta). \end{aligned} \quad (2)$$

We begin by defining minimal exponential families, a special class of exponential families with nice properties. To be minimal, the sufficient statistics must be linearly independent. We will later relax the requirement that we consider only minimal exponential families.

Definition 1. An exponential family of distributions generating a random variable $\mathbf{x} \in \mathcal{X}$ with $S(\mathbf{x}) \in \mathbb{R}^d$ is said to be minimal if $\exists \phi \in \mathbb{R}^d, \phi \neq 0$ s.t. $\exists c \in \mathbb{R}$ s.t. $\forall \mathbf{x} \in \mathcal{X} \phi^\top S(\mathbf{x}) = c$.

Next we present a few simple algebraic results of minimal exponential families.

Lemma 1. For two distributions p, q from the same minimal exponential family,

$$KL(p||q) = A(\theta_q) - A(\theta_p) - (\theta_q - \theta_p)^\top \nabla A(\theta_p) \quad (3)$$

where θ_p, θ_q are the natural parameters of p and q , and $A(\theta)$ is the log-normalizer for the exponential family.

Lemma 2. A minimal exponential family distribution satisfies these equalities:

$$\nabla A(\theta) = E_{Pr(\mathbf{x}|\theta)}[S(\mathbf{x})]$$

$$\nabla^2 A(\theta) = cov_{Pr(\mathbf{x}|\theta)}(S(\mathbf{x})).$$

Lemma 3. For a minimal exponential family distribution, its log-normalizer $A(\theta)$ is a strictly convex function over the natural parameters. This implies a bijection between θ and $E_{Pr(\mathbf{x}|\theta)}[S(\mathbf{x})]$.

These are standard results coming from some algebraic manipulations as seen in (Brown, 1986), and we omit the proof of these lemmas. Lemma 3 immediately leads to a useful corollary about minimal families and their conjugate prior families.

Corollary 4. For a minimal exponential family distribution, the conjugate prior family given in equation (2) is also minimal.

PROOF:

$T(\theta) = (\theta, -A(\theta))$ forms the sufficient statistics for the conjugate prior. Since $A(\theta)$ is strictly convex, there can be no linear relationship between the components of θ and $A(\theta)$. Definition 1 applies. \square

Our next result looks at sufficient conditions for getting a KL divergence of 0 in the limit when adding a finite perturbation vector γ to the natural parameters. The limit is taken over N , which will later be tied to the amount of data used in forming the posterior. As we now discuss posterior distributions also forming exponential families, our natural parameters will now be denoted by η and the random variables are now θ .

Lemma 5. Let $p(\theta|\eta)$ denote the distribution from an exponential family of natural parameter η , and let γ be a constant vector of the same dimensionality as η , and let η_N be

a sequence of natural parameters. If for every ζ on the line segment connecting η and $\eta + \gamma$ we have the spectral norm $\|\nabla^2 B(\zeta)\| < D_N$ for some constant D_N , then

$$KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N)) \leq D_N \|\gamma\|.$$

PROOF: This follows from noticing that equation (3) in Lemma 1 becomes the first-order Taylor approximation of $B(\eta_N)$ centered at $B(\eta_N + \gamma)$. From Taylor's theorem, there exists α between η_N and $\eta_N + \gamma$ such that $\frac{1}{2}\gamma^\top \nabla^2 B(\alpha)\gamma$ is equal to the error of this approximation.

$$B(\eta_N) = B(\eta_N + \gamma) + (-\gamma)^\top \nabla B(\eta_N + \gamma) + \frac{1}{2}\gamma^\top \nabla^2 B(\alpha)\gamma \quad (4)$$

From rearranging equation (3),

$$B(\eta_N + \gamma) = B(\eta_N) - KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N)) + (\gamma)^\top \nabla B(\eta_N + \gamma) \quad (5)$$

Using this substitution in (4) gives

$$B(\eta_N) = B(\eta_N) - KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N)) + \frac{1}{2}\gamma^\top \nabla^2 B(\alpha)\gamma. \quad (6)$$

Solving for $KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N))$ then gives the desired result:

$$KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N)) = \frac{1}{2}\gamma^\top \nabla^2 B(\alpha)\gamma \leq D_N \|\gamma\|.$$

□

This provides the heart of our results: If $\|\nabla^2 B(\zeta)\|$ is small for all ζ connecting η and $\eta + \gamma$, then we can conclude that $KL(p(\theta|\eta_N + \gamma)||p(\theta|\eta_N))$ is small with respect to $\|\gamma\|$. We wish to show that for η_N arising from observing N data points we have D_N approaching 0 as N grows. To achieve this, we will analyze a relationship between the norm of the natural parameter η and the covariance of the distribution it parameterizes. This relationship shows that posteriors with plenty of observed data have low covariance over $T(\theta)$, which permits us to use Lemma 5 to bound the KL divergence of our perturbed posteriors. Before we reach this relationship, first we prove that our posteriors have a well-defined mode, as our later relationship will require this mode to be well-behaved.

Lemma 6. Let $Pr(\mathbf{x}|\theta) = h(\mathbf{x}) \exp(\theta^\top S(\mathbf{x}) - A(\theta))$ be a likelihood function for θ and let there be a conjugate

prior $Pr(\theta|\eta) = f(\theta) \exp(\eta^\top T(\theta) - B(\eta))$, where both distributions are minimal exponential families. Let M be the space of natural parameters η , and Θ be the space of θ . Furthermore, assume η is the parameterization arising from the natural conjugate prior, such that $\eta = (\alpha\chi, \alpha)$. If the following conditions hold:

1. η is in the interior of M
2. $\alpha > 0$
3. $A(\theta)$ is a real, continuous, and differentiable
4. $B(\eta)$ exists, the distribution $Pr(\theta|\eta)$ is normalizable.

then

$$\operatorname{argmax}_{\theta \in \Theta} \eta^\top T(\theta) = \theta_\eta^*$$

is a well-defined function of η , and θ_η^* is in the interior of Θ .

PROOF:

Using our structure for the conjugate prior from (2), we can expand the expression $\eta^\top T(\theta)$.

$$\eta^\top T(\theta) = \alpha\chi^\top \theta - \alpha A(\theta)$$

We note that the first term is linear in θ , and that by minimality and Lemma 3, $A(\theta)$ is strictly convex. This implies $\eta^\top T(\theta)$ is strictly concave over θ . Thus any interior local maximum must also be the unique global maximum.

The gradient of with $\eta^\top T(\theta)$ respect to θ is simple to compute.

$$\nabla(\eta^\top T(\theta)) = \alpha\chi^\top - \alpha \nabla A(\theta)$$

This expression can be set to zero, and solving for θ_η^* shows it must satisfy

$$\nabla A(\theta_\eta^*) = \chi. \quad (7)$$

We remark by Lemma 2 that $\nabla A(\theta_\eta^*)$ is equal to $E_{Pr(\mathbf{x}|\theta_\eta^*)}[S(\mathbf{x})]$, and so this is the θ that generates a distribution with mean χ .

By the strict concavity, this is sufficient to prove θ_η^* is a unique local maximizer and thus the global maximum.

To see that θ_η^* must be in the interior of Θ , we use the fact that $A(\theta)$ is continuously differentiable. This means $\nabla A(\theta)$ is a continuous function of θ . Since η is in the interior of M , we can construct an open neighborhood around χ . The preimage of an open set under a continuous function

is also an open set, so this implies an open neighborhood exists around θ_η^* .

□

Now that we know θ_η^* is well defined for η in the interior of M , we can express our relationship on high magnitude posterior parameters and the covariance of the distribution over $T(\theta)$ they generate.

Lemma 7. *Let $Pr(\mathbf{x}|\theta) = h(\mathbf{x}) \exp(\theta^\top S(\mathbf{x}) - A(\theta))$ be a likelihood function for θ and let there be a conjugate prior $Pr(\theta|\eta) = f(\theta) \exp(\eta^\top T(\theta) - B(\eta))$, where both distributions are minimal exponential families. Let M be the space of natural parameters η , and Θ be the space of θ . Furthermore, assume η is the parameterization arising from the natural conjugate prior, such that $\eta = (\alpha\chi, \alpha)$.*

If $\exists \eta_0, \delta_1 > 0, \delta_2 > 0$ such that the conditions of Lemma 6 hold for $\eta \in \mathcal{B}(\eta_0, \delta_1)$, and we have these additional assumptions,

1. *the cone $\{k\eta' | k > 1, \eta' \in \overline{\mathcal{B}(\eta_0, \delta_1)}\}$ lies entirely in M*
2. *$A(\theta)$ is differentiable of all orders*
3. *$\exists P$ s.t. $\forall \theta \in \cup_{\eta' \in \overline{\mathcal{B}(\eta_0, \delta_1)}} \mathcal{B}(\theta_{\eta'}^*, \delta_2)$ all partial derivatives up to order 7 of $A(\theta)$ have magnitude bounded by P*
4. *$\exists w > 0$ such that $\forall \theta \in \cup_{\eta' \in \overline{\mathcal{B}(\eta_0, \delta_1)}} \mathcal{B}(\theta_{\eta'}^*, \delta_2)$ we have $\det(\nabla^2 A(\theta)) > w$*

then there exists C, K such that for $k > K$ the following bound holds $\forall \eta \in \mathcal{B}(\eta_0, \delta_1)$:

$$\|cov(T(\theta)|k\eta)\| < \frac{C}{k}.$$

PROOF:

This result follows from the Laplace approximation method for $B(\eta) = \int_{\Theta} e^{\eta^\top T(\theta)} d\theta$. The inner details of this approximation are show in Lemma 11. Here we show that our setting satisfies all the regularity assumptions for this approximation. First we define functions $s(\theta, \eta)$ and $F_k(\eta)$.

$$s(\theta, \eta) = \eta^\top T(\theta) = \alpha\chi^\top \theta - \alpha A(\theta) \quad (8)$$

$$\begin{aligned} F_k(\eta) &= B(k\eta) \\ &= \int_{\Theta} e^{k\eta^\top T(\theta)} d\theta \\ &= \int_{\Theta} e^{ks(\theta, \eta)} d\theta \end{aligned} \quad (9)$$

With these definitions, we may now begin to check the assumptions of Lemma 11 hold. We copy these assumptions

below, with a substitution of θ for ϕ and η for Y . The full details of Lemma 11 can be found at the end of section B.1.

1. $\phi_Y^* = \operatorname{argmax}_{\phi \in M} s(\phi, Y) = g(Y)$, a function of Y .
2. $\phi_{Y'}^*$ is in the interior of M for all $Y' \in \mathcal{B}(Y_0, \delta_1)$.
3. $g(Y)$ is continuously differentiable over the neighborhood $\mathcal{B}(Y_0, \delta_1)$.
4. $s(\phi, Y')$ has derivatives of all orders for $Y' \in \mathcal{B}(Y_0, \delta_1), \phi \in \mathcal{B}(\phi_{Y'}^*, \delta_2)$ and all partial derivatives up to order 7 are bounded by some constant P on this neighborhood.
5. $\exists w > 0$ such that $\forall Y' \in \mathcal{B}(Y_0, \delta_1), \forall \phi \in \mathcal{B}(\phi_{Y'}^*, \delta_2)$ we have $\det(\nabla_{\phi}^2 s(\phi, Y)) > w$.
6. $F_1(Y')$ exists for $Y' \in \mathcal{B}(Y_0, \delta_1)$, the integral is finite.

We now show these conditions hold one-by-one. Let η denote an arbitrary element of $B(\eta_0, \delta)$.

1. θ_η^* is a well-defined function (Lemma 6).
2. θ_η^* is in the interior of Θ (Lemma 6).
3. $g(\eta)$ follows the inverse of $\nabla A(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}^d$. This vector mapping has a Jacobian $\nabla^2 A(\theta)$ which assumption 4 guarantees has non-zero determinant on this neighborhood. This satisfies the Inverse Function Theorem to show $g(\eta)$ is continuously differentiable.
4. $s(\theta, \eta)$ has derivatives of all orders, and are suitably bounded as s is composed of a linear term and the differentiable function $A(\theta)$, where we have bounded the derivatives of $A(\theta)$.
5. Assumption 4 from this lemma translates directly.
6. $F_1(\eta) = B(\eta)$ which exists by virtue of η being in the space of valid natural parameters.

This completes all the requirements of Lemma 11, which guarantees the existence of C and K such that for any $k > K$ and any $\eta \in \mathcal{B}(\eta_0, \delta_1)$, if we let ψ denote $k\eta$, we have:

$$\|\nabla_{\psi}^2 B(\psi)\| = \|\nabla_{\psi}^2 \log F_k(\psi/k)\| < \frac{C}{k}.$$

We conclude by noting that $\nabla_{\psi}^2 B(\psi)$ is the covariance of the posterior with parameterization $\psi = k\eta$.

□

Now that all our machinery is in place, it remains to be seen under what conditions the posterior satisfies the conditions of the previous Lemmas, along with extending to the case where γ is a random variable, and not just a fixed finite vector.

Lemma 8. For a minimal exponential family given a conjugate prior, where the posterior takes the form $Pr(\theta|\mathbf{X}, \chi, \alpha) \propto g(\theta)^{n+\alpha} \exp\left(\theta^\top\left(\sum_{i=1}^n S(\mathbf{x}^{(i)}) + \alpha\chi\right)\right)$, where $p(\theta|\eta)$ denotes this posterior with a natural parameter vector η , if there exists a $\delta > 0$ such that these assumptions are met:

1. the data \mathbf{X} comes i.i.d. from a minimal exponential family distribution with natural parameter $\theta_0 \in \Theta$
2. θ_0 is in the interior of Θ
3. the function $A(\theta)$ has all derivatives for θ in the interior of Θ
4. $cov_{Pr(\mathbf{x}|\theta)}(S(\mathbf{x}))$ is finite for $\theta \in \mathcal{B}(\theta_0, \delta)$
5. $\exists w > 0$ s.t. $\det(cov_{Pr(\mathbf{x}|\theta)}(S(\mathbf{x}))) > w$ for $\theta \in \mathcal{B}(\theta_0, \delta)$
6. the prior $Pr(\theta|\chi, \alpha)$ is integrable and has support on a neighborhood of θ^*

then for any mechanism generating a perturbed posterior $\tilde{p}_N = p(\theta|\eta_N + \gamma)$ against a noiseless posterior $p_N = p(\theta|\eta_N)$ where γ comes from a distribution that does not depend on the number of data observations N and has finite covariance, this limit holds:

$$\lim_{N \rightarrow \infty} E[KL(\tilde{p}_N || p_N)] = 0.$$

PROOF:

We begin by fixing the randomness of the noise γ that the mechanism will add to the natural parameters of the posterior.

We wish to show that the KL divergence goes to zero in the limit, which we will achieve by showing that for large enough data sizes, both the perturbed and unperturbed posteriors lie w.h.p. in a region where we can use Lemmas 5 and 7 apply.

To compute the posterior, after drawing a collection \mathbf{X} of N data observations, we compute the sum of the sufficient statistics and add them to the prior's parameters.

$$\eta_N = \left(\alpha\chi + \sum S(\mathbf{x}^{(i)}), \alpha + N\right)$$

η_N is a random variable depending on the data observations \mathbf{X} . To analyze how it behaves, a couple related random variables will be defined, all implicitly conditioned on the constant θ_0 . Let \mathbf{Y} denote a random variable matching the distribution of a single observation, and let $\mathbf{U}_N = \frac{1}{N} \sum S(\mathbf{x}^{(i)})$ which has covariance $\frac{1}{N} cov(S(\mathbf{Y}))$. The expected value for \mathbf{U}_N is of course $E[S(\mathbf{Y})]$.

By a vector version of the Chebyshev inequality for a random vector \mathbf{U} , (Chen, 2007)

$$Pr\left(\left(\mathbf{U} - E[\mathbf{U}]\right)^\top (cov(\mathbf{U}))^{-1} (\mathbf{U} - E[\mathbf{U}]) \geq \nu\right), \leq \frac{d}{\nu}, \quad (10)$$

where d is the dimensionality of \mathbf{U} . Using the spectral norm $\|(cov(\mathbf{U}_N))^{-1}\|$ and the l_2 norm $\|\mathbf{U}_N - E[\mathbf{U}_N]\|$ with some rearrangement, we can show the following inequalities. We note that the covariance of \mathbf{U}_N must be invertible, since the covariance of \mathbf{Y} is invertible by assumption (5).

$$Pr\left(\|\mathbf{U}_N - E[\mathbf{U}_N]\| \cdot \|(cov(\mathbf{U}_N))^{-1}\| \geq \nu\right) \leq \frac{d}{\nu} \quad (11)$$

$$Pr\left(\|\mathbf{U}_N - E[\mathbf{U}_N]\| \geq \nu \|cov(\mathbf{U}_N)\| \right) \leq \frac{d}{\nu} \quad (12)$$

$$Pr\left(\|\mathbf{U}_N - E[S(\mathbf{Y})]\| \geq \frac{\nu}{N} \|cov(\mathbf{Y})\| \right) \leq \frac{d}{\nu} \quad (13)$$

Thus for any $\epsilon > 0, \tau > 0$, there exists $N_{\epsilon, \tau}$ such that when the number of data observations N exceeds $N_{\epsilon, \tau}$

$$Pr(\|\mathbf{U}_N - E[\mathbf{Y}]\| \geq \epsilon) \leq \tau. \quad (14)$$

We now define two modified vectors of natural parameters $\eta_a = \frac{\eta_N}{N} = (\mathbf{U}_N, 1) + \frac{1}{N}(\alpha\chi, \alpha)$ and $\eta_b = \frac{\eta_N + \gamma}{N} = (\mathbf{U}_N, 1) + \frac{1}{N}(\alpha\chi, \alpha) + \frac{1}{N}\gamma$. From these definitions, one can see

$$E[\eta_a] = (E[\mathbf{Y}], 1) + \frac{1}{N}(\alpha\chi, \alpha)$$

$$E[\eta_b] = E[\eta_a] + \frac{1}{N}\gamma$$

$$\|\eta_a - (E[\mathbf{Y}], 1)\| \leq \|(\mathbf{U}_N, 1) - (E[\mathbf{Y}], 1)\| + \frac{1}{N}\|\alpha\chi\| \quad (15)$$

$$\|\eta_b - (E[\mathbf{Y}], 1)\| \leq \|(\mathbf{U}_N, 1) - (E[\mathbf{Y}], 1)\| + \frac{1}{N}(\|\alpha\chi\| + \|\gamma\|). \quad (16)$$

From the concentration bound in (14), we know η_a and η_b can be made to lie w.h.p. in a region near their expectations with large N , and we wish to show this region satisfies all

the regularity assumptions seen in Lemma 7. Lemma 6 states θ_η^* is a continuously differentiable function of η . Let it be denoted by the function $r(\eta)$. For $\eta_0 = (E[\mathbf{Y}], 1)$, we see from equation (7) that $r(\eta_0) = \theta_0$.

The preimage $r^{-1}(\mathcal{B}(\theta_0, \delta))$ is an open set, since it is the continuous preimage of an open set. Thus there exists δ' such that $\mathcal{B}(\eta_0, \delta') \subset r^{-1}(\mathcal{B}(\theta_0, \delta/2))$.

We may now pick $\epsilon \leq \delta'/2$ and let $N'_{\delta', \tau} = \max(\frac{2}{\delta'}(\|\gamma\| + \|\alpha\chi\|), N_{\epsilon, \tau})$. When $n > N'_{\delta', \tau}$, we have $\frac{1}{N}\|\alpha\chi\| + \frac{1}{N}\|\gamma\| \leq \delta'/2$ and (14), (15), (16) together show the following:

$$Pr(\eta_a \notin \mathcal{B}(\eta_0, \delta') \vee \eta_b \notin \mathcal{B}(\eta_0, \delta')) \leq \tau. \quad (17)$$

With high probability, η_a and η_b both lie in a neighborhood of η_0 . Further, all η in this neighborhood have modes $\theta_\eta^* \in \mathcal{B}(\theta_0, \delta)$, a region that assumptions (4) and (5) tell us is well-behaved. The assignment $\delta_1 = \delta'$ and $\delta_2 = \delta/2$ satisfies the conditions for Lemma 7 with assumptions (2),(3),(4),(5),(6) serving to round out the rest of the regularity assumptions of Lemma 7 with trivial translations.

By the construction, we have $\eta_N = N\eta_a$ and $\eta_N + \gamma = N\eta_b$. For any ζ on the line segment connecting η_N and $\eta_N + \gamma$, we have $\zeta = N\eta_c$ for some η_c on the line segment connecting η_a and η_b .

Therefore by Lemma 7, there exists a K and a C such that if $N > K$ we have $\|cov(T(\theta)|\zeta)\| < \frac{C}{N}$. This bound can be used in Lemma 5 with $D_N = O(1/N)$ to see

$$KL(\tilde{p}_N||p_N) = O(1/N)C\|\gamma\|$$

whenever $N > \max(N'_{\delta', \tau}, K)$ with arbitrarily high probability $1 - \tau$. Letting τ approach 0, we can extend this to the expectation over the randomness of \mathbf{X} , as with probability 1 our random variables will lie in the region where this inequality holds.

$$\limsup_{N \rightarrow \infty} E_{\mathbf{X}}[KL(\tilde{p}_N||p_N)] = 0 \quad (18)$$

Equation (18) is w.r.t. to a fixed γ , but the desired result is an expectation over γ and \mathbf{X} . First, let us express this expectation in terms of γ and \mathbf{X} . Letting $D_N = O(1/N)$ denote the bound used in Lemma 5 and N being sufficiently large:

$$\begin{aligned} E[KL(\tilde{p}_N||p_N)] &= \int E_{\mathbf{X}}[KL(\tilde{p}_N||p_N)|\gamma] dPr(\gamma) \\ &\leq \int D_N \|\gamma\| dPr(\gamma). \end{aligned} \quad (19)$$

The assumption that γ comes from a distribution of finite variance ensures the right side of (19) is integrable. By an application of Fatou's Lemma, the following inequality holds:

$$\begin{aligned} &\int \limsup_{N \rightarrow \infty} E_{\mathbf{X}}[KL(\tilde{p}_N||p_N)|\gamma] dPr(\gamma) \\ &\geq \limsup_{N \rightarrow \infty} \int E_{\mathbf{X}}[KL(\tilde{p}_N||p_N)|\gamma] dPr(\gamma). \end{aligned} \quad (20)$$

The left hand side has been shown to be zero by equations (18) and (19), and the right hand side is bounded below by 0 since KL divergences are never negative. Thus this inequality suffices to show the limit is zero and prove the desired result.

□

Corollary 9. *The Laplace mechanism on an exponential family satisfies the noise distribution requirements of Lemma 8 when the sensitivity of the sufficient statistics is finite and either the exponential family is minimal, or if the exponential family parameters θ are identifiable.*

PROOF: If the exponential family is already minimal, this result is trivial. If it is not minimal, there exists a minimal parameterization. We wish to show adding noise to the non-minimal parameters is equivalent to adding differently distributed noise to the minimal parameterization, and this new noise distribution also satisfies the noise distribution requirements of Lemma 8: the noise distribution does not depend on N and it has finite covariance.

Let us explicitly construct a minimal parameterization for this family of distributions. If the exponential family is not minimal, this means the d dimensions of the sufficient statistics $S(\mathbf{x})$ of the data are not fully linearly independent. Let $S(x)_j$ be the j^{th} component of $S(\mathbf{x})$ and k be the maximal number of linearly independent sufficient statistics, and without loss of generality assume they are the first k components. Let $\tilde{S}(\mathbf{x})$ be the vector of these k linearly independent components.

For $\forall j > k, \forall x \exists \phi_j \in \mathbb{R}^k$ such that $S(x)_j = \phi_j \cdot \tilde{S}(\mathbf{x}) + z_j$. We wish to build a minimal exponential family distribution that is identical to the original one, but is parameterized only by $\tilde{S}(\mathbf{x})$ as the sufficient statistics and some $\tilde{\theta}$ as the natural parameters. For these two distributions to be equivalent for all x , it suffices to have equality on the exponents.

$$(\theta^\top S(\mathbf{x}) - A(\theta)) = (\tilde{\theta}^\top \tilde{S}(\mathbf{x}) - \tilde{A}(\tilde{\theta})) \quad (21)$$

Examining the difference of the two sides, we get

$$\begin{aligned}
& \theta^\top S(x) - \tilde{\theta}^\top \tilde{S}(x) - A(\theta) + \tilde{A}(\tilde{\theta}) \\
&= \sum_{j=1}^k (\theta_j - \tilde{\theta}_j) S(x)_j + \sum_{j=k+1}^d \theta_j S(x)_j - A(\theta) + \tilde{A}(\tilde{\theta}).
\end{aligned} \tag{22}$$

Using the known linear dependence for $j > k$, this can be rewritten as

$$\begin{aligned}
& \sum_{j=1}^k (\theta_j - \tilde{\theta}_j) S(\mathbf{x})_j + \sum_{j=k+1}^d \theta_j (\phi_j \cdot \tilde{S}(\mathbf{x}) + z_j) \\
& \qquad \qquad \qquad - A(\theta) + \tilde{A}(\tilde{\theta})
\end{aligned} \tag{23}$$

$$\begin{aligned}
&= \sum_{j=1}^k (\theta_j - \tilde{\theta}_j) S(\mathbf{x})_j + \sum_{j=k+1}^d \theta_j (\phi_j \cdot \tilde{S}(\mathbf{x})) \\
& \qquad \qquad \qquad + \sum_{j=k+1}^d \theta_j z_j - A(\theta) + \tilde{A}(\tilde{\theta}).
\end{aligned} \tag{24}$$

Now since $\tilde{S}(\mathbf{x})$ is merely the first k components of $S(\mathbf{x})$, the first two sums of (24) are each simply dot products of $\tilde{S}(\mathbf{x})$ and can be combined as $(\theta_{[k]} - \tilde{\theta} + \sum_{j=k+1}^d \theta_j \phi_j)^\top \tilde{S}(\mathbf{x})$ where $\theta_{[k]}$ is the vector of the first k components of θ . We can force equation (21) to hold by choosing $\tilde{\theta}$ and $\tilde{A}(\tilde{\theta})$ appropriately to set equation (24) to zero.

$$\begin{aligned}
\tilde{\theta} &= \theta_{[k]} + \sum_{j=k+1}^d \theta_j \phi_j \\
\tilde{A}(\tilde{\theta}) &= - \sum_{j=k+1}^d \theta_j z_j + A(\theta)
\end{aligned}$$

We note that this requires $\tilde{A}(\tilde{\theta})$ to truly be a function depending only on $\tilde{\theta}$, but we have written it in terms of θ instead. This is justifiable by the assumption that the natural parameters θ are identifiable, that is each distribution over \mathbf{x} is associated with just one $\theta \in \Theta$. This means there is a bijection from θ and $\tilde{\theta}$, which ensures $\tilde{A}(\tilde{\theta})$ is a well-defined function.

This suffices to characterize the way the additional natural parameters affect the parameters of the equivalent minimal system. Any additive noise to a component θ_j translates linearly to additive noise on the components $\tilde{\theta}_j$, meaning the Laplace mechanism's noise distribution on the non-minimal parameter space still corresponds to some noise distribution on the minimal parameters that does not depend on the data size N , and it still has a finite covariance. If the minimal exponential family tends towards a KL divergence of zero, the equivalent non-minimal exponential family must as well. \square

Theorem B.1. *Under the assumptions of Lemma 8, the Laplace mechanism has an asymptotic posterior of $\mathcal{N}(\theta_0, 2\mathbb{I}^{-1}/N)$ from which drawing a single sample has an asymptotic relative efficiency of 2 in estimating θ_0 , where \mathbb{I} is the Fisher information at θ_0 .*

PROOF:

The assumptions of Lemma 8 match the Laplace regularity assumptions under which asymptotic normality holds, and we know that the unperturbed posterior p_N converges to $\mathcal{N}(\theta^*, 2\mathbb{I}^{-1}/N)$ under the Bernstein-von Mises theorem (Kass et al., 1990). If \tilde{p}_N is the posterior of the Laplace mechanism for a fixed randomness, then we have $\lim_{N \rightarrow \infty} KL(\tilde{p}_N || p_N) = 0$ and \tilde{p}_N must converge to the same distribution as p_N . From this it is clear that samples from p_N and from \tilde{p}_N both have an asymptotic relative efficiency of 2. We once again argue that if this asymptotic behavior holds for any fixed randomness of the Laplace mechanism, it also holds for the Laplace mechanism as a whole. \square

To show the previous results, we relied on some mathematical results involving the covariances of posteriors after observing a large amount of data. We still need to show these bounds on the covariances, which will be accomplished by adapting existing Laplace approximation methods. Before we get there, we will need one quick result about convex functions with a positive definite Hessian in order to perform the approximation:

Lemma 10. *Let $f(y) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a strictly convex function with minimum at y^* . If $\nabla^2 f(y^*)$ is positive definite and $\nabla^3 f(y)$ exists everywhere, then for any $c > 0$ there exists $b > 0$ such that $|f(y) - f(y^*)| \leq b$ implies $\|y - y^*\| \leq c$.*

PROOF:

By the existence of $\nabla^3 f(y)$ and thus the continuity of $\nabla^2 f(y)$, we know there exists a positive $\delta < c$ and a $w > 0$ such that $y \in B(y^*, \delta)$ implies $\nabla^2 f(y) - w\mathbb{I}$ is positive semi-definite, where \mathbb{I} is the identity matrix. (i.e. the spectral norm $\|\nabla^2 f(y)\| \geq w$)

As y^* is the global minimum, we know the gradient is 0 at y^* . Thus for $y \in B(y^*, \delta)$ this leads to a Taylor expansion of the form

$$\begin{aligned}
f(y) &= f(y^*) + (y - y^*)^\top \nabla f(y') + \frac{1}{2} (y - y^*)^\top \nabla^2 f(y') (y - y^*) \\
&\geq f(y^*) + \frac{w}{2} \|y - y^*\|^2
\end{aligned} \tag{25}$$

for some y' on the line segment connecting y and y^* . The inequality follows from the second derivative being positive definite on this neighborhood.

Consider the set $Q_\epsilon = \{y \text{ s.t. } \|y - y^*\| = \epsilon\}$. By equation (25) we know for $y \in Q_\epsilon$ we have $|f(y) - f(y^*)| \geq \frac{w\epsilon}{2}$ if $\epsilon \leq \delta$.

For any $y \notin B(y^*, \delta)$, there exists $t \in (0, 1)$ such that $(1-t)y^* + ty \in Q_\delta$ by the continuity of the norm.

By strict convexity, we know

$$tf(y) + (1-t)f(y^*) > f(ty + (1-t)y^*)$$

$$f(y) > \frac{1}{t}f(ty + (1-t)y^*) + \frac{t-1}{t}f(y^*)$$

$$f(y) - f(y^*) > \frac{1}{t}f(ty + (1-t)y^*) - \frac{1}{t}f(y^*).$$

If we let t satisfy $(1-t)y^* + ty \in Q_\delta$ we know $t = \delta/\|y - y^*\| \leq 1$. Substituting with (25) we get

$$f(y) - f(y^*) > \frac{(w/2)\delta + f(y^*)}{t} - \frac{1}{t}f(y^*) = \frac{w\delta}{2t} \geq \frac{w\delta}{2}.$$

Thus if we let $b = \frac{w\delta}{2}$, we see $\|y - y^*\| > c$ implies $|f(y) - f(y^*)| > b$.

The desired result then follows as the contrapositive.

□

Lemma 10 will be used to demonstrate a regularity assumption required in the next lemma, which performs all the heavy lifting in using the Laplace approximation. Lemma 11 adapts a previous argument about Laplace approximations of a posterior. This adapted Laplace approximation argument forms the core of Lemma 7, which allows us to see the covariance of posteriors shrink as more data is observed.

Lemma 11. *Let $s(\phi, Y)$ be a function $M \times U \rightarrow \mathbb{R}$, where M is the space of ϕ and U is the space of Y .*

For functions of the form $F_k(Y) = \int_{\phi \in M} e^{ks(\phi, Y)} d\phi$, if the following regularity assumptions hold for some $\delta_1 > 0$, $\delta_2 > 0$, $Y_0 \in M$:

1. $\phi_Y^* = \operatorname{argmax}_{\phi \in M} s(\phi, Y) = g(Y)$, a function of Y
2. ϕ_Y^* , is in the interior of M for all $Y' \in \mathcal{B}(Y_0, \delta_1)$
3. $g(Y)$ is continuously differentiable over the neighborhood $\mathcal{B}(Y_0, \delta_1)$
4. $s(\phi, Y')$ has derivatives of all orders for $Y' \in \mathcal{B}(Y_0, \delta_1)$, $\phi \in \mathcal{B}(\phi_Y^*, \delta_2)$ and all partial derivatives up to order 7 are bounded by some constant P on this neighborhood

5. $\exists w > 0$ such that $\forall Y' \in \mathcal{B}(Y_0, \delta_1), \forall \phi \in \mathcal{B}(\phi_Y^*, \delta_2)$ we have $\det(\nabla_\phi^2 s(\phi, Y)) > w$

6. $F_1(Y')$ exists for $Y' \in \mathcal{B}(Y_0, \delta_1)$, the integral is finite

then there exists C and K such that for any $k > K$ and any $Y' \in \mathcal{B}(Y_0, \delta_1)$, letting $\psi = kY'$, the spectral norm $\|\nabla_\psi^2 \log F_k(\psi/k)\| < \frac{C}{k}$.

PROOF:

Our goal here is to bound $\|\nabla_\psi^2 \log F_k(\psi/k)\|$, which we will achieve by characterizing $F_k(\psi/k)$ and analyzing its derivatives.

We will be using standard Laplace approximation methods seen in (Kass et al., 1990) to explore $F_k(\psi)$. To begin, we must show our assumptions satisfy the regularity assumptions for the approximation.

For a fixed $Y' \in \mathcal{B}(Y_0, \delta)$, from condition 5 we know there exists a neighborhood around ϕ_Y^* where $\nabla_\phi^2 s(\phi, Y)$ is positive definite. For $\delta' > 0$, let $Q_{\delta', Y} = \{\phi \in M \text{ s.t. } \|\phi - \phi_Y^*\| \leq \delta'\}$. By using Lemma 10 we can verify the following expression for any $\delta' \in (0, \delta)$:

$$\limsup_{k \rightarrow \infty} \sup_{\phi \notin Q_{\delta', Y}} s(\phi, Y) - s(\phi_Y^*, Y) < 0. \quad (26)$$

Note that the right hand side does not depend on k , and Lemma 10 guarantees a non-zero bound for the right hand side for any $\delta' \in (0, \delta)$. Equation (26) exactly matches condition (iii)' of Kass, and its intuitive meaning is that for any δ' , there exists sufficiently large k such that the integral F_k is negligible outside the region $Q_{\delta'}$.

Conditions (4),(5),(6) also match directly the conditions given by Kass, though we note we require even higher derivatives to be bounded or present. These extra derivatives will be used later to extend the argument given by Kass to suit our purposes and give a uniform bound across a neighborhood.

Theorem 1 of (Kass et al., 1990) gives the following result, when we set their b to the constant 1:

$$F_k(Y) = (2\pi)^{\frac{m}{2}} [\det(k\nabla^2 s(\phi_Y^*, Y))]^{-\frac{1}{2}} \exp(-ks(\phi_Y^*, Y)) Z(kY) \quad (27)$$

$$Z(kY) = 1 + \frac{1}{k} \left(\frac{1}{72} \sum (\nabla_\phi^3 s(\phi_Y^*, Y))_{(pqr)} (\nabla^3 s(\phi_Y^*, Y))_{(def)} \mu_{pqrdef}^6 - \frac{1}{24} \sum (\nabla^4 s(\phi_Y^*, Y))_{(defg)} \mu_{defg}^4 \right) + O(k^{-2}), \quad (28)$$

where m is the dimensionality of Y , μ_{pqrd}^6 and μ_{def}^4 are the sixth and fourth central moments of a multivariate Gaussian with covariance matrix $(\nabla^2 s(\phi_Y^*, Y))^{-1}$. All sums are written in the Einstein summation notation. We remark that the $O(k^{-2})$ error term of this approximation also depends on kY .

What we are really interested in is the quantity $\nabla_\psi^2 \log F_k(\psi)$ evaluated at $\psi = kY$. We take the logarithm of (27):

$$\begin{aligned} \log F_k(\psi/k) &= \log \left((2\pi)^{\frac{m}{2}} [\det(k\nabla^2 s(\phi_Y^*, Y))]^{-\frac{1}{2}} \right. \\ &\quad \left. \cdot \exp(-ks(\phi_Y^*, Y))Z(\psi) \right) \\ &= \log \left((2\pi)^{\frac{m}{2}} \right) - \frac{1}{2} \log([\det(k\nabla^2 s(\phi_Y^*, Y))]) \\ &\quad - ks(\phi_Y^*, Y) + \log(Z(\psi)). \end{aligned} \quad (29)$$

We define new functions $\tilde{s}_0, \tilde{s}_1, \tilde{s}_2$ to simplify the analysis.

$$\tilde{s}_0(Y) = s(\phi_Y^*, Y) = s(g(Y), Y) \quad (30)$$

$$\tilde{s}_1(Y) = \nabla_\phi s(\phi_Y^*, Y) = \nabla_\phi s(g(Y), Y) \quad (31)$$

$$\tilde{s}_2(Y) = \nabla_\phi^2 s(\phi_Y^*, Y) = \nabla_\phi^2 s(g(Y), Y) \quad (32)$$

By assumptions (3) and (4) we know these functions are continuously differentiable on $\mathcal{B}(Y_0, \delta_1)$ as they are the composition of continuously differentiable functions on the compact set $\mathcal{B}(Y_0, \delta_1)$.

We next look at the first derivative of (29). We remark that the partial derivatives of $\log \det(X)$ are given by $X^{-\top}$.

$$\begin{aligned} \nabla_\psi \log F_k(\psi/k) &= \nabla_\psi \left[-\frac{1}{2} \log([\det(k\tilde{s}_2(\psi/k))]) \right. \\ &\quad \left. - \nabla_\psi [k\tilde{s}_0(\psi/k)] + \nabla_\psi \log(Z(\psi)) \right] \\ &= -\frac{1}{2} (k\tilde{s}_2(\psi/k))^{-\top} \frac{1}{k} \\ &\quad + \tilde{s}_1(\psi/k) + \frac{\nabla_\psi(Z(\psi))}{Z(\psi)} \end{aligned} \quad (33)$$

Now that we have an expression for $\nabla_\psi \log F_k(\psi/k)$, we take yet another derivative w.r.t. to ψ to get our desired ∇_ψ^2 .

$$\begin{aligned} \nabla_\psi^2 \log F_k(\psi/k) &= \nabla_\psi \left[-\frac{1}{2} (k\tilde{s}_2(\psi/k))^{-\top} \frac{1}{k} \right. \\ &\quad \left. + \nabla_\psi [\tilde{s}_1(\psi/k)] \right. \\ &\quad \left. + \nabla_\psi \left[\frac{\nabla_\psi(Z(\psi))}{Z(\psi)} \right] \right] \end{aligned} \quad (34)$$

Let us consider each of the three terms on the right side of (34) in isolation. For the first term, we introduce yet another function $\tilde{s}_{-2}(Y)$, the composition of \tilde{s}_2 with the matrix inversion.

$$\tilde{s}_{-2}(Y) = (\tilde{s}_2(Y))^{-1}$$

With this new function in hand, we further condense the first term of (34).

$$\begin{aligned} \nabla_\psi \left[-\frac{1}{2} (k\tilde{s}_2(\psi/k))^{-\top} \frac{1}{k} \right] &= \nabla_\psi \left[-\frac{1}{2k} (\tilde{s}_{-2}(\psi/k)) \frac{1}{k} \right] \\ &= -\frac{1}{2k^3} \nabla_Y \tilde{s}_{-2}(\psi/k) \\ &= O(k^{-3}) \end{aligned} \quad (35)$$

We previously remarked that \tilde{s}_2 is continuously differentiable on the compact set $\mathcal{B}(Y_0, \delta_1)$. Condition (5) informs us that $\tilde{s}_2(Y)$ is bounded away from being a singular matrix on $\mathcal{B}(Y_0, \delta_1)$, so the matrix inversion is also uniformly continuous on this compact set. This means $\nabla_Y \tilde{s}_{-2}(\psi/k)$ has a finite supremum over $\mathcal{B}(Y_0, \delta_1)$ and thus we can say this term is $O(k^{-3})$ uniformly on this neighborhood.

Next we consider the second term of (34).

$$\nabla_\psi [\tilde{s}_1(\psi/k)] = \frac{1}{k} \tilde{s}_2(\psi/k) = O(k^{-1}) \quad (36)$$

From the continuity of $\tilde{s}_2(\psi/k)$ on our compact neighborhood, we know $\tilde{s}_2(Y)$ has a finite supremum over the compact set $\mathcal{B}(Y_0, \delta_1)$, which gives the uniform $O(k^{-1})$ bound.

Finally, we must consider the third term of (34).

$$\nabla_\psi \left[\frac{\nabla_\psi(Z(\psi))}{Z(\psi)} \right] = \frac{\nabla^2(Z(\psi))}{Z(\psi)} - \frac{\nabla(Z(\psi))(\nabla(Z(\psi)))^\top}{Z(\psi)^2} \quad (37)$$

Recall that $Z(\psi)$ had a local $O(k^{-2})$ error term as given by (Kass et al., 1990). We wish to bound the derivatives of $\log F_k(\psi)$, but the local bound on this error term given by Kass does not bound its derivatives. However, a slight modification of the argument of (Kass et al., 1990) shows that our added assumptions about the higher order derivatives are sufficient to control the behavior of this error term. The following expression is their equation (2.2), translated to our setting:

$$\begin{aligned}
& \exp(-ks(\phi, Y)) = \\
& \exp(-ks(\phi_Y^*, Y)) \exp\left(\frac{1}{2}\nabla^2 s(\phi_Y^*, Y)u^2\right)W(\phi, Y) \quad (38) \\
& W(\phi, Y) = 1 - \frac{1}{6}k^{-1/2}\nabla^3 s(\phi_Y^*, Y)u^3 \\
& \quad + \frac{1}{72}k^{-1}(\nabla^3 s(\phi_Y^*, Y))^2u^6 \\
& \quad - \frac{1}{24}k^{-1}\nabla^4 s(\phi_Y^*, Y)u^4 \\
& \quad - \frac{1}{120}k^{-3/2}\nabla^5 s(\phi_Y^*, Y)u^5 \\
& \quad + \frac{1}{72}k^{-3/2}\nabla^3 A(s(\phi_Y^*, Y))\nabla^4 s(\phi_Y^*, Y)u^7 \\
& \quad + G(\phi, \phi_Y^*, Y), \quad (39)
\end{aligned}$$

where $G(\phi, \phi_Y^*, Y)$ is the fifth-order Taylor expansion error term (i.e. it depends on the sixth-order partial derivatives at some ϕ' between ϕ and ϕ_Y^*).

We may continue this Taylor expansion another degree further to bound the variation of $G(\phi, \phi_Y^*, Y)$ for $\phi \in \mathcal{B}(\phi_Y^*, \delta_2)$. We will consider $Z(\psi)$, $\nabla_\psi Z(\psi)$, and $\nabla_\psi^2 Z(\psi)$ as three separate functions, each permitting a higher order Taylor expansion. Each will have their own respective error term depending on the seventh-order partial derivatives at some ϕ' , but we note that ϕ' is not necessarily the same for each of them.

The argument of (Kass et al., 1990) already shows how the terms composing their $O(k^{-2})$ error term can be bounded in terms of $\nabla_\phi^6 S(\phi_Y^*, Y)$. It is trivial to show an analogous result for our higher order approximations. This allows us to extend our approximation of $Z(\psi)$ and its derivatives uniformly to the neighborhood $\mathcal{B}(\phi_Y^*, \delta_2)$. The newly introduced extra approximation terms are $O(k^{-v})$ with $v \geq 2$, and so our uniform bounds are still simply $O(k^{-2})$, though with a larger constant now.

Let k be sufficiently large, and let Q, R, S be positive constants satisfying $0 < Q < \|Z(\psi)\|$, $R > k\|\nabla_\psi Z(\psi)\|$, $S > k\|\nabla_\psi^2 Z(\psi)\|$ for all ψ in $\{\psi|\psi/k \in B(Y_0, \delta)\}$. We remark that Q exists by virtue of $Z = 1 + O(k^{-1}) + O(k^{-2})$. R and S similarly exist by $\|\nabla_\psi Z(\psi)\|$ and $\|\nabla_\psi^2 Z(\psi)\|$ both being $O(k^{-1})$ with no constant term in front.

$$\nabla_\psi \left[\frac{\nabla_\psi(Z(\psi))}{Z(\psi)} \right] \leq \frac{S}{kQ} - \frac{R^2}{k^2Q^2} \text{ for all } Y' \in B(Y_0, \delta)$$

This right hand side is clearly $O(k^{-1})$, and we have uniform bounds across our neighborhood.

$$\nabla_\psi \left[\frac{\nabla_\psi(Z(\psi))}{Z(\psi)} \right] = O(k^{-1}) \quad (40)$$

Combining the results of (35), (36), (40) with their sum in (34), we get this result:

$$\|\nabla_\psi^2 \log F_k(\psi/k)\| = O(k^{-1}). \quad (41)$$

This uniform asymptotic bound then ensures we have the intended result: $\exists C, K$ such that $\forall Y \in \mathcal{B}(Y_0, \delta_1)$ when $k > K$ and $\psi = kY$ we have $\|\nabla_\psi^2 \log F_k(\psi/k)\| \leq C/k$

□

C PRIVACY PROPERTIES OF OTHER MCMC ALGORITHMS

In the main manuscript we showed the privacy cost of Gibbs sampling by interpreting it as an instance of the exponential mechanism. Here, we show the privacy cost of two other widely used MCMC algorithms: Metropolis-Hastings and annealed importance sampling.

C.1 METROPOLIS-HASTINGS UPDATES

Since Gibbs updates are a special case of Metropolis-Hastings updates, one might conjecture that general Metropolis-Hastings updates may be differentially private as well. However, the accept/reject decision contains a subtle non-determinacy which violates pure- ϵ differential privacy. Consider a Metropolis-Hastings update with a symmetric proposal $\theta' \sim f(\theta, \theta')$ (a.k.a. a Metropolis update),

$$Pr(\text{accept}; \mathbf{X}, \theta, \theta', T) = \min\left(1, \left(\frac{Pr(\theta'|\mathbf{X})}{Pr(\theta|\mathbf{X})}\right)^{\frac{1}{T}}\right) \quad (42)$$

where T is the temperature of the Markov chain. For these updates, ‘‘uphill’’ moves are never rejected. Since a move may be uphill in one database and downhill in a neighbor, we cannot bound the ratio of reject decisions, which violates differential privacy. It turns out that Metropolis updates do have a weaker privacy guarantee, by resorting to (ϵ, δ) -differential privacy:

Theorem C.1. *Let \mathbf{X} be private data and θ be a public current value of the variables we wish to infer. A Metropolis update invariant to the posterior $Pr(\theta|\mathbf{X})$ at temperature $T = \frac{2\Delta \log Pr(\theta, \mathbf{X})}{\epsilon}$, with symmetric proposal $\theta' \sim f(\theta, \theta')$ and with $Pr(\text{reject}; \mathbf{X}, \theta, T) = \int f(\theta, \theta')(1 - Pr(\text{accept}; \mathbf{X}, \theta, \theta', T))d\theta' \leq \delta$, is (ϵ, δ) -differentially private.*

A proof of Theorem C.1 is provided below in Appendix D. Essentially, we can bound the ratio of probabilities for accept decisions under neighboring databases, but not for reject decisions. If rejections are rare, these privacy-violating outcomes are rare, which is sufficient for (ϵ, δ) -privacy. On the other hand, δ must be very small for a meaningful level of privacy, e.g. less than the inverse of any polynomial in

the number of data points N (Dwork and Roth, 2013), so this may not typically correspond to a practical privacy-preserving sampling algorithm.

C.2 ANNEALED IMPORTANCE SAMPLING

The privacy results for Gibbs sampling and Metropolis-Hastings updates reveal a close connection between privacy and the temperature of the Markov chain. Low-temperature chains are high-fidelity but privacy-expensive, while high-temperature chains are low-fidelity but privacy-cheap, and also mix more rapidly. This suggests that annealing methods, such as annealed importance sampling (AIS) (Neal, 2001), may be effective in this context, by allowing savings in the privacy budget in the early iterations of MCMC while also traversing the state space more rapidly. AIS is a Monte Carlo method which anneals from a high-temperature distribution to the target distribution (in our case the posterior) via MCMC updates at a sequence of temperatures, producing importance weights for each sample to correct for the annealing. AIS takes as input an annealing path, a sequence of unnormalized distributions $f_n(\theta), \dots, f_0(\theta)$ at different temperatures. We can obtain a privacy-preserving AIS annealing path by varying ϵ :

$$f_j(\theta) = Pr(\theta, \mathbf{X})^{\beta_j}, \beta_j = \frac{\epsilon_j}{2\Delta \log Pr(\theta, \mathbf{X})}, \quad (43)$$

where each intermediate distribution f_j is an instance of Equation 6, and ϵ_j is the privacy cost for an exact sample from f_j . We can sample at each temperature using the private Gibbs transition operator from Equation 19. The privacy cost of an AIS sample is computed via the composition theorem,

$$\epsilon^{(AIS)} = \sum_j \sum_l \epsilon_j = \sum_j D\epsilon_j, \quad (44)$$

where l ranges over the D variables to be updated. If each Gibbs update only depends on a single data point \mathbf{x}_l , we can improve this via parallel composition (Song et al., 2013) to

$$\epsilon^{(AIS)} = \sum_j \epsilon_j. \quad (45)$$

On completion of the algorithm we must compute importance weights ω_i for the samples $\theta^{(i)}$:

$$\begin{aligned} \log \omega_i &= \sum_{j=0}^{n-1} \left(\log f_j(\theta^{(i,j)}) - \log f_{j+1}(\theta^{(i,j)}) \right) \quad (46) \\ &= \sum_{j=0}^{n-1} \left(\beta_j \log Pr(\theta^{(i,j)}, \mathbf{X}) - \beta_{j+1} \log Pr(\theta^{(i,j)}, \mathbf{X}) \right) \\ &= \frac{1}{2\Delta \log Pr(\theta, \mathbf{X})} \sum_{j=0}^{n-1} (\epsilon_j - \epsilon_{j+1}) \log Pr(\theta^{(i,j)}, \mathbf{X}). \end{aligned}$$

We only need to release private copies of the importance weights at the end of the procedure, as they are not used during the algorithm. If we are not interested in computing normalization constants, we can release a normalized version of the weights, dividing by $\sum_i \omega_i$. This is a discrete distribution which sums to one, and so it lives on the simplex. This has $L1$ sensitivity at most 2, and can be protected by the Laplace mechanism. Another possible alternative is to perform resampling of the $\theta^{(i)}$'s according to this distribution, approximated and protected via the exponential mechanism.

D PROOF OF THEOREM C.1

Here, we prove the differential privacy result for Metropolis-Hastings given in Theorem C.1, above. PROOF: Let $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}$ be neighboring databases. By the definition of differential privacy, we need to bound the ratios of the probability of each outcome, $\{(\text{accept, reject}), \theta^{(new)}\}$ for these two databases. We consider accept and reject outcomes separately.

D.1 ACCEPT OUTCOME

The probability of an accepted move to location $\theta^{(new)} = \mathbf{z}'$ is

$$\begin{aligned} Pr(\text{accept}, \theta^{(new)} = \theta'; \mathbf{X}, \theta) \\ = f(\theta, \theta') Pr(\text{accept}; \mathbf{X}, \theta, \theta', T). \end{aligned}$$

We must bound the probability ratio of this outcome under the two neighboring datasets. Consider first a slightly simpler question, the ratio of probabilities for an accept decision, having already selected the proposal θ' ,

$$\frac{Pr(\text{accept}; \mathbf{X}^{(1)}, \theta, \theta')}{Pr(\text{accept}; \mathbf{X}^{(2)}, \theta, \theta')}.$$

We will perform the computation in log space. We have the log of the acceptance probabilities as

$$\begin{aligned} \log Pr(\text{accept}; \mathbf{X}, \theta, \theta', T) &= \\ \min \left(0, \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} (\log Pr(\theta'|\mathbf{X}) - \log Pr(\theta|\mathbf{X})) \right) \\ &= \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \min \left(0, \log Pr(\theta'|\mathbf{X}) - \log Pr(\theta|\mathbf{X}) \right). \end{aligned}$$

The difference in log probabilities for the accept outcome is

$$\begin{aligned} \log Pr(\text{accept}; \mathbf{X}^{(1)}, \theta, \theta') - \log Pr(\text{accept}; \mathbf{X}^{(2)}, \theta, \theta') \\ = \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \times \\ \left(\min \left(0, \log Pr(\theta'|\mathbf{X}^{(1)}) - \log Pr(\theta|\mathbf{X}^{(1)}) \right) \right. \\ \left. - \min \left(0, \log Pr(\theta'|\mathbf{X}^{(2)}) - \log Pr(\theta|\mathbf{X}^{(2)}) \right) \right). \end{aligned}$$

Let

$$\begin{aligned} a &= \log Pr(\theta'|\mathbf{X}^{(1)}) - \log Pr(\theta|\mathbf{X}^{(1)}) \\ b &= \log Pr(\theta'|\mathbf{X}^{(2)}) - \log Pr(\theta|\mathbf{X}^{(2)}) . \end{aligned}$$

There are four cases to consider:

$$a \leq 0, b \leq 0:$$

$$\min(0, a) - \min(0, b) = a - b$$

$$a > 0, b \leq 0:$$

$$\min(0, a) - \min(0, b) = -b \leq a - b$$

$$a \leq 0, b > 0:$$

$$\min(0, a) - \min(0, b) = a \leq 0$$

$$a > 0, b > 0:$$

$$\min(0, a) - \min(0, b) = 0 .$$

So either $\min(0, a) - \min(0, b) \leq 0$, in which case the difference in log probabilities is $\leq 0 \leq \epsilon$, or $\min(0, a) - \min(0, b) \leq a - b$. In the former, we are done, so consider the latter case:

$$\begin{aligned} & \log Pr(\text{accept}; \mathbf{X}^{(1)}, \theta, \theta') - \log Pr(\text{accept}; \mathbf{X}^{(2)}, \theta, \theta') \\ & \leq \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \left((\log Pr(\theta'|\mathbf{X}^{(1)}) - \log Pr(\theta|\mathbf{X}^{(1)})) \right. \\ & \quad \left. - (\log Pr(\theta'|\mathbf{X}^{(2)}) - \log Pr(\theta|\mathbf{X}^{(2)})) \right) \\ & = \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \left(\log Pr(\theta'|\mathbf{X}^{(1)}) - \log Pr(\theta'|\mathbf{X}^{(2)}) \right. \\ & \quad \left. + \log Pr(\theta|\mathbf{X}^{(2)}) - \log Pr(\theta|\mathbf{X}^{(1)}) \right) \\ & = \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \left(\log Pr(\theta'|\mathbf{X}^{(1)}) - \log Pr(\theta'|\mathbf{X}^{(2)}) \right. \\ & \quad \left. + \log Pr(\mathbf{X}^{(1)}) - \log Pr(\mathbf{X}^{(2)}) \right. \\ & \quad \left. + \log Pr(\theta|\mathbf{X}^{(2)}) - \log Pr(\theta|\mathbf{X}^{(1)}) \right. \\ & \quad \left. + \log Pr(\mathbf{X}^{(2)}) - \log Pr(\mathbf{X}^{(1)}) \right) \\ & = \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \left(\log Pr(\theta', \mathbf{X}^{(1)}) - \log Pr(\theta', \mathbf{X}^{(2)}) \right. \\ & \quad \left. + \log Pr(\theta, \mathbf{X}^{(2)}) - \log Pr(\theta, \mathbf{X}^{(1)}) \right) \\ & \leq \frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})} \left(\Delta \log Pr(\theta, \mathbf{X}) + \Delta \log Pr(\theta, \mathbf{X}) \right) \\ & = \epsilon . \end{aligned}$$

The inequality in the last line follows from Equation 7 in the main paper.

Having bounded the log ratio of probabilities by ϵ for the simpler case where the proposal θ' is given, we can

now bound the ratios for the full output, of the form $(\text{accept}, \theta^{(new)})$, as required for ϵ -differential privacy, by simply cancelling the log transition probabilities:

$$\begin{aligned} & \log Pr(\text{accept}, \theta^{(new)} = \theta'; \mathbf{X}^{(1)}, \theta) \\ & - \log Pr(\text{accept}, \theta^{(new)} = \theta'; \mathbf{X}^{(2)}, \theta) \\ & = \log f(\theta, \theta') + \log Pr(\text{accept}; \mathbf{X}^{(1)}, \theta, \theta') \\ & - (\log f(\theta, \theta') + \log Pr(\text{accept}; \mathbf{X}^{(2)}, \theta, \theta')) \\ & = \log Pr(\text{accept}; \mathbf{X}^{(1)}, \theta, \theta') - \log Pr(\text{accept}; \mathbf{X}^{(2)}, \theta, \theta') \\ & \leq \epsilon . \end{aligned}$$

This is as desired for pure- ϵ privacy, and so the weaker (ϵ, δ) -criterion holds for this outcome as well.

D.2 REJECT OUTCOME

If we could also similarly bound the difference in log probabilities between neighboring databases for the outcome $(\text{reject}, \theta^{(new)} = \theta)$ by ϵ , then the Metropolis update would be ϵ -differentially private. Consider first the reject probabilities after the proposal θ' is selected:

$$\begin{aligned} Pr(\text{reject}; \mathbf{X}, \theta, \theta') & = 1 - \min \left(1, \left(\frac{Pr(\theta'|\mathbf{X})}{Pr(\theta|\mathbf{X})} \right)^{\frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})}} \right) \\ & = \max \left(0, 1 - \left(\frac{Pr(\theta'|\mathbf{X})}{Pr(\theta|\mathbf{X})} \right)^{\frac{\epsilon}{2\Delta \log Pr(\theta, \mathbf{X})}} \right) . \end{aligned}$$

When $Pr(\theta'|\mathbf{X}) > Pr(\theta|\mathbf{X})$, the probability of a reject decision is 0. It is possible to construct scenarios where the probability of a reject decision is 0 for all proposals θ' , e.g. when θ is at a global minimum, so we cannot in general lower bound the overall probability of a reject,

$$\begin{aligned} Pr(\text{reject}, \theta^{(new)} = \theta; \mathbf{X}, \theta) & = \int f(\theta, \theta') (1 - Pr(\text{accept}; \mathbf{X}, \theta, \theta', T)) d\theta' . \end{aligned}$$

If $Pr(\text{reject}, \theta^{(new)} = \theta; \mathbf{X}, \theta) = 0$ occurs in database $\mathbf{X}^{(1)}$ and not in $\mathbf{X}^{(2)}$, the ratio of probabilities for this outcome will be infinite due to a division by 0, violating ϵ -differential privacy. Under our assumptions, we have an additional guarantee that $Pr(\text{reject}, \theta^{(new)} = \theta; \mathbf{X}, \theta) \leq \delta$, i.e. the probability of a rejection outcome, and therefore the probability of an outcome that violates ϵ -differential privacy, is less than δ . To demonstrate (ϵ, δ) privacy and complete the proof, we observe that this condition implies that the (ϵ, δ) -criterion holds for the reject outcome:

$$\begin{aligned} Pr(\text{reject}, \theta^{(new)} = \theta; \mathbf{X}^{(1)}, \theta) & \leq \delta \leq \exp(\epsilon) Pr(\text{reject}, \theta^{(new)} = \theta; \mathbf{X}^{(2)}, \theta) + \delta . \end{aligned}$$

□

$x_{i,d}^{(r,t)}$	Discrete-valued feature d of log entry i , from region r , timestep t .
$z_{r,t}$	Latent state at region r , timestep t .
$A_{k,k'}$	Transition probability from state k to k' .
$\theta_j^{(k,d)}$	Discrete emission probability for cluster k 's d 'th feature being outcome j .
α, β	Dirichlet concentration parameters.
$N_{r,t}$	Number of log entries (observations) in region r at timestep t .
D	Number of features in the observations.
K	Number of latent clusters.

Table 1: Notation for the Wikileaks naive Bayes HMM model.

E DETAILS OF WIKILEAKS WAR LOGS HMM

In this appendix we describe the technical details of the HMM model with naive Bayes observations, which we apply to the Wikileaks War Logs data. The assumed generative process of the model is:

For $k = 1, \dots, K$ //For each latent cluster

$\mathbf{A}_{k,:} \sim \text{Dirichlet}(\alpha)$ // K -dimensional

For $d = 1, \dots, D$ //For each feature

$\theta^{(k,d)} \sim \text{Dirichlet}(\beta)$ // K_d -dimensional

$\mathbf{A}_{0,:} \sim \text{Dirichlet}(\alpha)$ //Dummy state

For $r = 1, \dots, R$ //For each region

$z_{r,0} = 0$ //Dummy initial state

For $t = 1, \dots, T$ //For each timestep

$z_{r,t} \sim \text{Discrete}(\mathbf{A}_{z_{r,t-1},:})$

For $i = 1, \dots, N_{r,t}$ //For each log entry

For $d = 1, \dots, D$ //For each feature

$x_{i,d}^{(r,t)} \sim \text{Discrete}(\theta^{(z_{r,t},d)})$.

Here, α and β correspond to the concentration parameters for appropriately dimensioned Dirichlet distributions. See Table 1 for a summary of the notation. The generative model corresponds to the joint probability

$$\begin{aligned}
Pr(\mathbf{A}, \theta, \mathbf{Z}, \mathbf{X} | \alpha, \beta) &= \\
&\prod_{k=0}^K Pr(\mathbf{A}_{k,:} | \alpha) \prod_{k=1}^K \prod_{d=1}^D Pr(\theta^{(k,d)} | \beta) \\
&\times \prod_{r=1}^R \prod_{t=1}^T Pr(z_{r,t} | z_{r,t-1}, \mathbf{A}) \\
&\times \prod_{r=1}^R \prod_{t=1}^T \prod_{i=1}^{N_{r,t}} Pr(x_{i,d}^{(r,t)} | z_{r,t}, \theta).
\end{aligned} \tag{47}$$

Inspired by Goldwater and Griffiths (2007), we marginalize out the transition matrix \mathbf{A} . Let $\mathbf{X}^{(r,t)}$ be an $N_{r,t} \times D$ matrix containing the log entry observations at region r ,

timestep t . We obtain the following partially collapsed Gibbs update for $z_{r,t}$:

$$\begin{aligned}
Pr(z_{r,t} | z_{r,t-1}, z_{r,t+1}, \mathbf{X}^{(r,t)}, \theta, \alpha) & \\
&\propto Pr(z_{r,t} | z_{r,t-1}) Pr(z_{r,t+1} | z_{r,t}) Pr(\mathbf{X}^{(r,t)} | z_{r,t}, \theta) \\
&= \frac{n_{z_{r,t}, z_{r,t-1}} + \alpha}{n_{z_{r,t-1}} + K\alpha} \frac{n_{z_{r,t+1}, z_{r,t}} + \alpha}{n_{z_{r,t}} + \mathbb{I}[z_{r,t-1} = z_{r,t} = z_{r,t+1}] + \alpha} \\
&\quad \times Pr(\mathbf{X}^{(r,t)} | z_{r,t}, \theta),
\end{aligned} \tag{48}$$

where $n_{z,z'}$ are transition counts, excluding the current z to be updated, and the transition probabilities are implicitly conditioned on all other z 's, which they depend on via the transition counts. The indicator functions arise from book-keeping as the counts are modified by changing the current state. Due to conjugacy we have a simple update for $\theta^{(k,d)}$,

$$Pr(\theta^{(k,d)} | \mathbf{X}, \mathbf{Z}, \beta) \sim \text{Dirichlet}(n_{d,k,:} + \beta), \tag{49}$$

where $n_{d,k,:} = \sum_{r,t} n_{r,t,d,:}$ is a K_d -dimensional count vector of counts for feature d in cluster k .

E.1 PRESERVING PRIVACY

To privatize the likelihood via the Laplace mechanism, we first write Equation 50 in exponential family form. The conditional likelihood for $\mathbf{X}^{(r,t)}$ given $z_{r,t}$ can be written as

$$\begin{aligned}
Pr(\mathbf{X}^{(r,t)} | z_{r,t}, \theta) &= \prod_{i=1}^{N_{r,t}} Pr(x_{i,d}^{(r,t)} | z_{r,t}, \theta) \\
&= \prod_{i=1}^{N_{r,t}} \prod_{d=1}^D \theta_{x_{i,d}^{(r,t)}}^{(z_{r,t},d)} \\
&= \prod_{d=1}^D \prod_{j=1}^{K_d} \theta_j^{(z_{r,t},d)^{n_{r,t,d,j}}},
\end{aligned} \tag{50}$$

where $n_{r,t,d,j} = \sum_{i=1}^{N_{r,t}} \mathbb{I}[x_{i,d}^{(r,t)} = j]$, and $\mathbb{I}[\cdot]$ is the indicator function. From here we obtain the exponential family form

$$Pr(\mathbf{X}^{(r,t)} | z_{r,t}, \theta) = \exp \left(\sum_{d=1}^D \sum_{j=1}^{K_d} n_{r,t,d,j} \log \theta_j^{(z_{r,t},d)} \right). \tag{51}$$

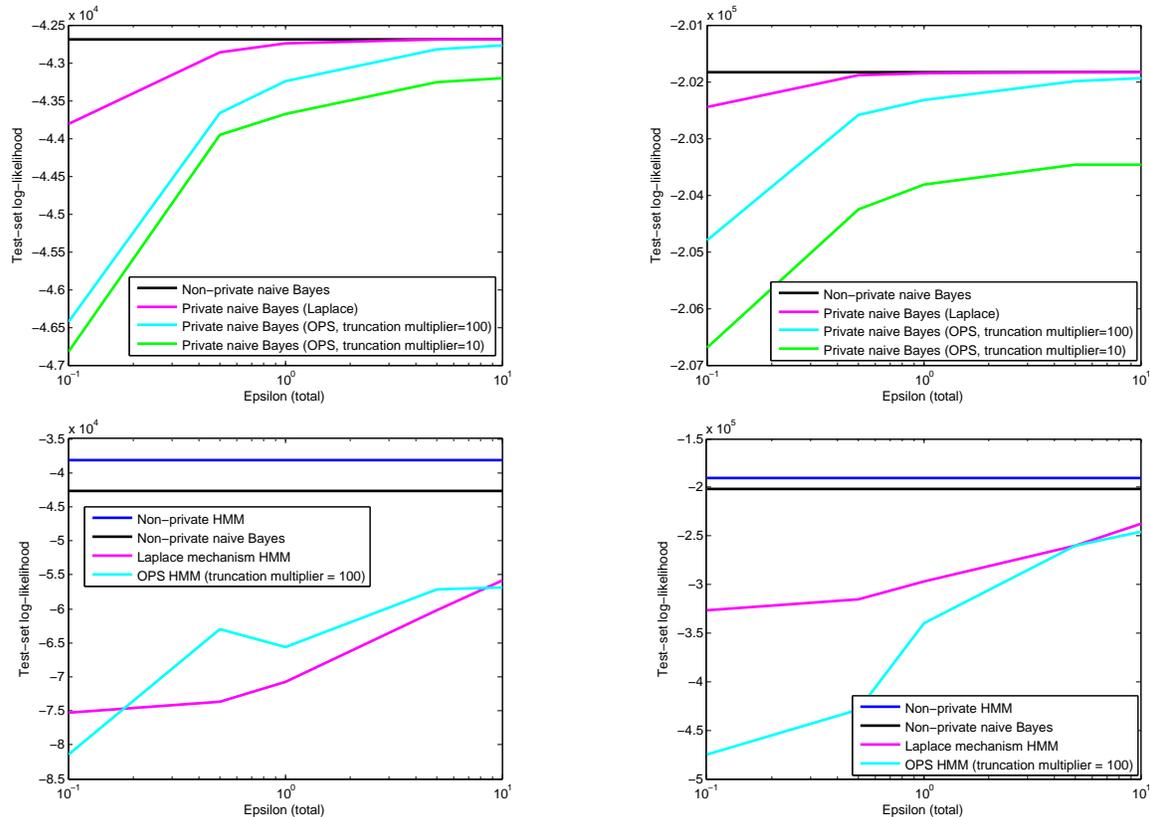


Figure 2: Log-likelihood of held-out data for a naive Bayes model, equivalent to the HMM with one timestep (**Top**) and the full HMM (**Bottom**). **Left:** Afghanistan. **Right:** Iraq. Truncation point for the truncated Dirichlet distributions for OPS was set to $a_0 = \frac{1}{MK_d}$, with truncation multiplier $M = 10$ and $M = 100$.

The sufficient statistics are the counts $n_{r,t,d,j}$, which we can privatize via the Laplace mechanism, resulting in private counts $\hat{n}_{r,t,d,j}$. As a sum of indicator vectors, each count vector $n_{r,t,d,:}$ has L1 sensitivity = 2. We can perform the Gibbs updates for \mathbf{Z} in a privacy-preserving manner by substituting the private counts for the counts in Equation 48. To preserve privacy when updating θ , Equation 49 can be estimated based on the privacy-preserving counts $\hat{n}_{r,t,d,:}$. Importantly, we only need to compute private counts $\hat{n}_{r,t,d,j}$ once, at the beginning of the algorithm, and these privatized counts can be reused for all of the Gibbs updates.

E.2 EXPERIMENTAL DETAILS

We performed some simple preprocessing steps before the experiment. Casualty count fields for each log entry were binarized (0 versus > 0). The wounded/killed/detained fields were merged disjunctively into one casualty indicator field. The *Friendly* (i.e. U.S. military) and *HostNation* (Iraq or Afghanistan) casualty indicators were combined into one field via disjunction. For the Iraq dataset, there were some missing data issues that had to be addressed. No data was available for the 5th

month, which was removed. Most regions had no data for the final year of the Iraq data, so this was also removed. Finally, we removed the MND-S and MND-NE region codes from our analysis, as these regions had very little data.

To simulate from truncated Dirichlet distributions for the Gibbs updates of the OPS method, we used the approach of Fang et al. (2000), which involves sequentially drawing each component based on a truncated Beta distribution. Full visualization results are shown in Figures 3 to 6. Log-likelihood results on held-out data are given in Figure 2. In this experiment, we randomly held-out 10% of the region/timestep pairs for testing for each of 5 train/test splits, and reported the average log-likelihood over the repeats.

References

Brown, L. D. (1986). Fundamentals of statistical exponential families with applications in statistical decision theory. *Lecture Notes-Monograph Series*, 9:i-279.

Chen, X. (2007). A new generalization of Chebyshev inequality for random vectors. *arXiv preprint arXiv:0707.0805*.

- Dwork, C. and Roth, A. (2013). The algorithmic foundations of differential privacy. *Theoretical Computer Science*, 9(3-4):211–407.
- Fang, K.-T., Geng, Z., and Tian, G.-L. (2000). Statistical inference for truncated Dirichlet distribution and its application in misclassification. *Biometrical journal*, 42(8):1053–1068.
- Goldwater, S. and Griffiths, T. (2007). A fully Bayesian approach to unsupervised part-of-speech tagging. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics (ACL)*, pages 744–751.
- Kass, R., Tierney, L., and Kadane, J. (1990). The validity of posterior expansions based on laplaces method. *Bayesian and likelihood methods in statistics and econometrics*, 7:473.
- Neal, R. M. (2001). Annealed importance sampling. *Statistics and Computing*, 11(2):125–139.
- Song, S., Chaudhuri, K., and Sarwate, A. D. (2013). Stochastic gradient descent with differentially private updates. In *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*, pages 245–248. IEEE.
- Wang, Y.-X., Fienberg, S. E., and Smola, A. (2015). Privacy for free: Posterior sampling and stochastic gradient Monte Carlo. *Proceedings of The 32nd International Conference on Machine Learning (ICML)*, pages 2493–2502.

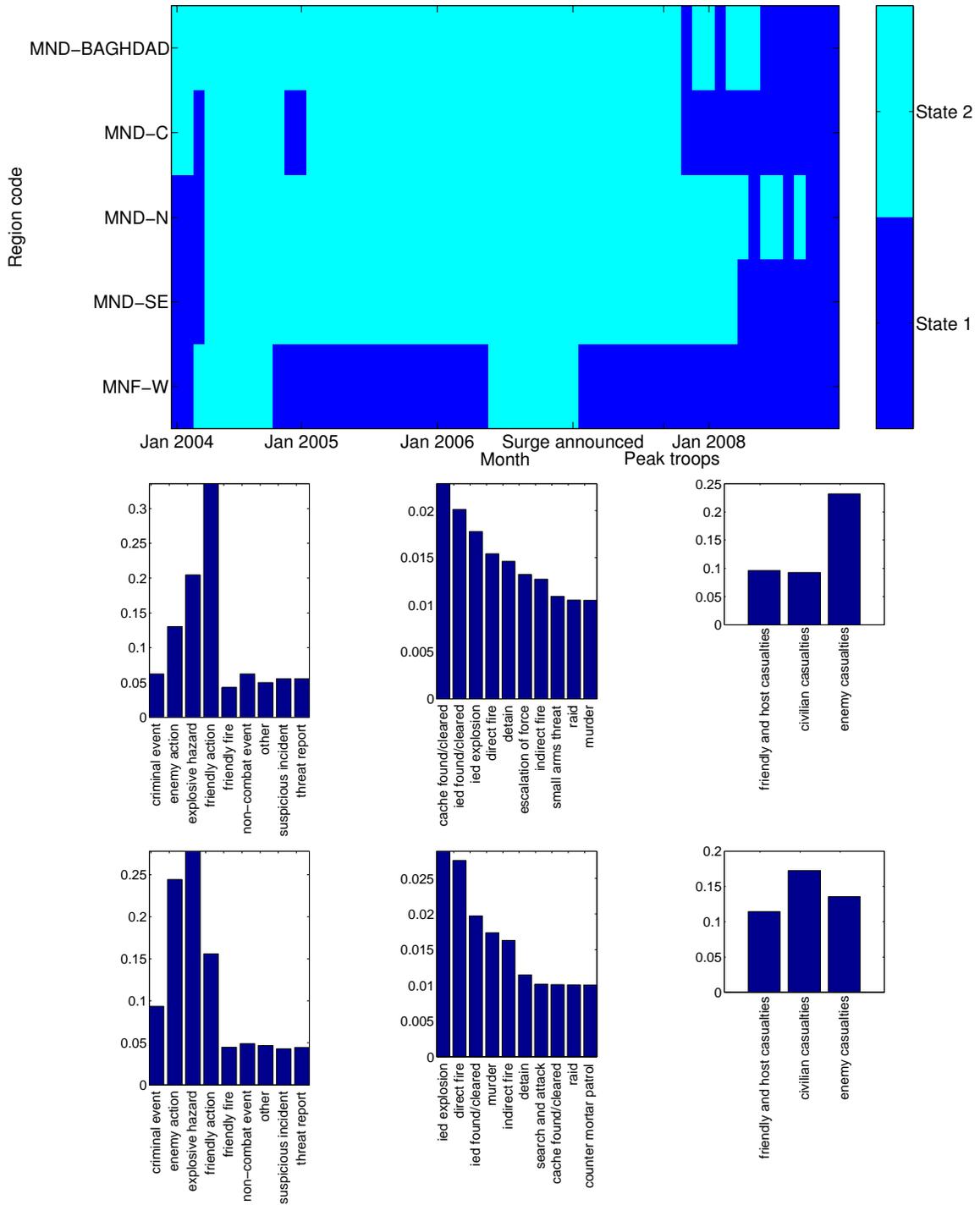


Figure 3: State assignments of privacy-preserving HMM on Iraq (Laplace mechanism, $\epsilon = 5$) (Top). Middle: State 1. Bottom: State 2.

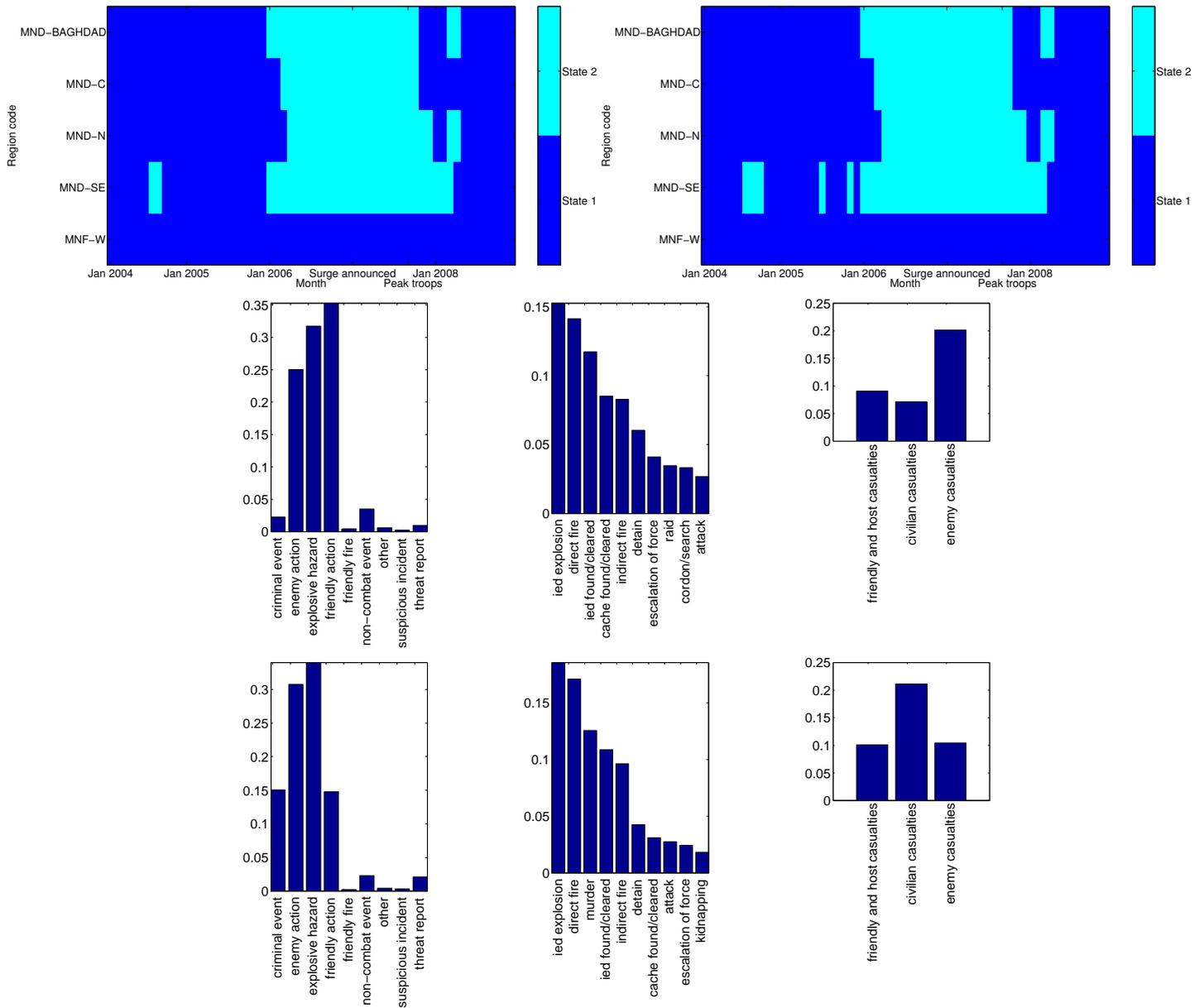


Figure 4: State assignments and parameters for OPS privacy-preserving HMM on Iraq. (OPS, $\epsilon = 5$, truncation point $a_0 = \frac{1}{100K_d}$). **Top Left:** Estimate from last 100 samples. **Top Right:** Estimate from last one sample. **Middle:** State 1. **Bottom:** State 2.

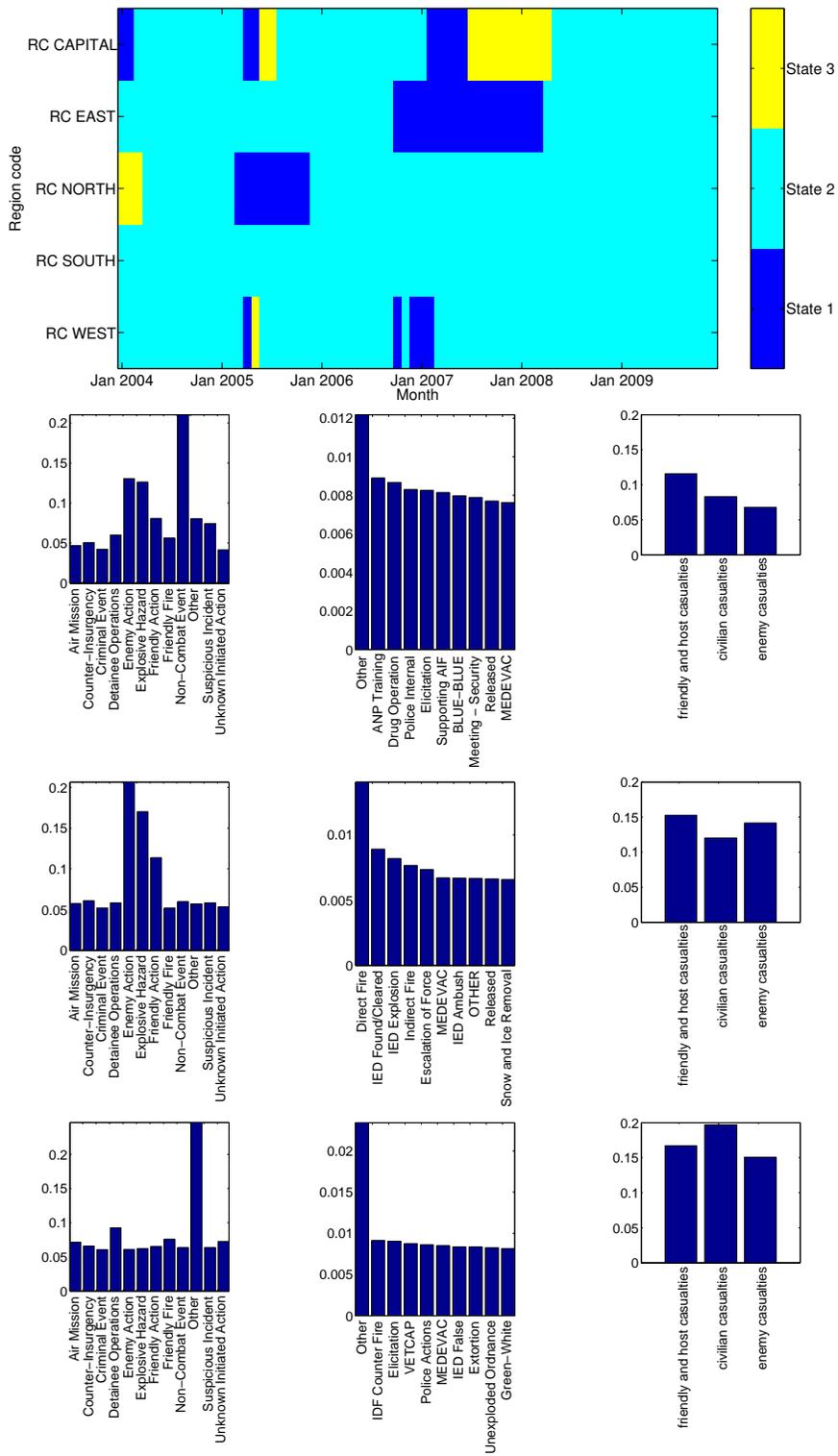


Figure 5: State assignments of privacy-preserving HMM on Afghanistan (Laplace mechanism, $\epsilon = 5$) (**Top**). Parameters for States 1, 2, and 3, ordered from top to bottom.

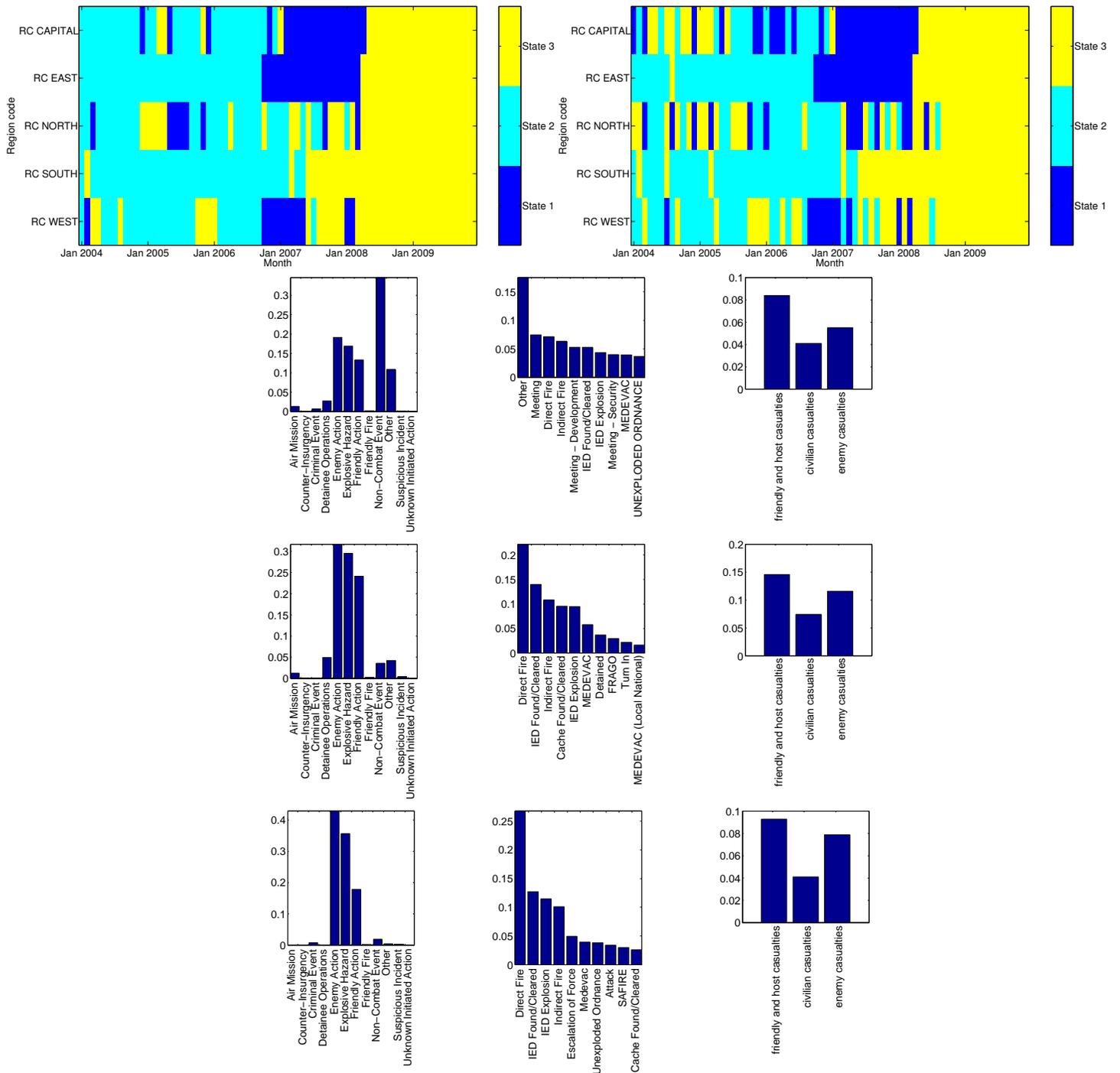


Figure 6: State assignments and parameters for OPS privacy-preserving HMM on Afghanistan. (OPS, $\epsilon = 5$, truncation point $a_0 = \frac{1}{100K_d}$). **Top Left:** Estimate from last 100 samples. **Top Right:** Estimate from last one sample. Parameters for States 1, 2, and 3, ordered from top to bottom.