

Variational Bayesian Unlearning

Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Patrick Jaillet

Presented by Ananth Mahadevan

February 10, 2021

Contents

- 1 Learning
 - VI and ELBO
- 2 Unlearning
 - Exact
 - EUBO
 - Adjusted EUBO
 - Reverse KL
- 3 Results

Bayesian Basics

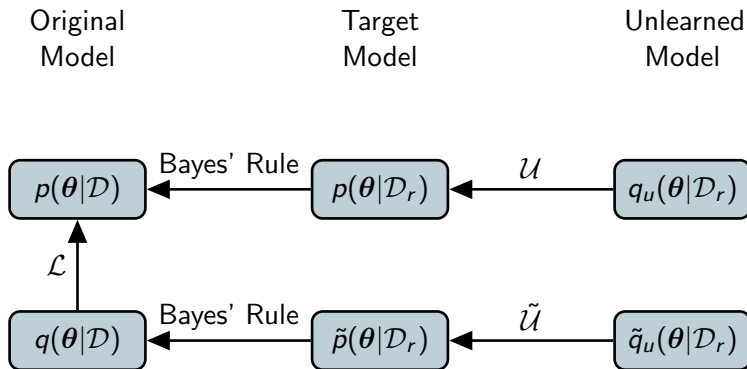
Learning

- Unknown model parameters θ
- Prior belief $p(\theta)$
- Set \mathcal{D} of training data
- Learn approximate posterior belief $q(\theta|\mathcal{D}) \approx p(\theta|\mathcal{D})$

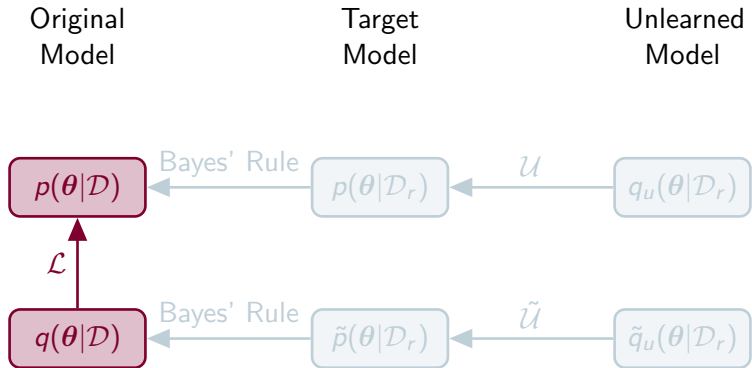
Unlearning

- \mathcal{D} partitioned into \mathcal{D}_e *erased data* and \mathcal{D}_r of *remaining data*
- $\mathcal{D} = \mathcal{D}_r \cup \mathcal{D}_e$ and $\mathcal{D}_r \cap \mathcal{D}_e = \emptyset$.
- Approximate $p(\theta|\mathcal{D}_r)$

Overview



Overview



Evidence Lower Bound (ELBO)

- Need to **minimize** the KL divergence

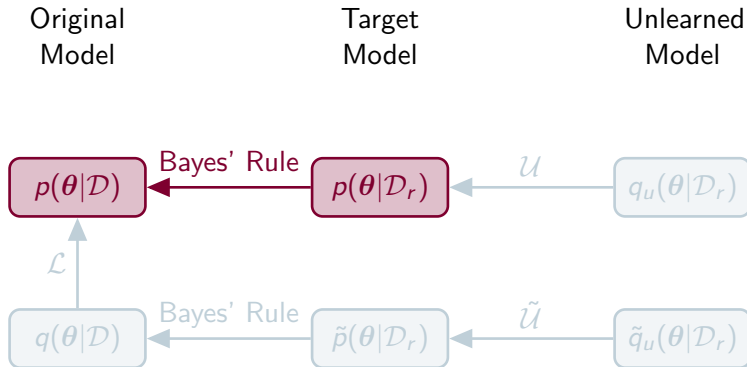
$$\text{KL}[q(\boldsymbol{\theta}|\mathcal{D}) \parallel p(\boldsymbol{\theta}|\mathcal{D})] \triangleq \int q(\boldsymbol{\theta}|\mathcal{D}) \log(q(\boldsymbol{\theta}|\mathcal{D})/p(\boldsymbol{\theta}|\mathcal{D})) \, d\boldsymbol{\theta}$$

- Or, **maximize** the *evidence lower bound* (ELBO)

$$\mathcal{L} \triangleq \underbrace{\int q(\boldsymbol{\theta}|\mathcal{D}) \log p(\mathcal{D}|\boldsymbol{\theta}) \, d\boldsymbol{\theta}}_{\text{increase likelihood}} - \underbrace{\text{KL}[q(\boldsymbol{\theta}|\mathcal{D}) \parallel p(\boldsymbol{\theta})]}_{\text{remember prior}}$$

- Use **Black Box Variational Inference** (BBVI) if \mathcal{L} cannot be evaluated in closed form
 - Stochastic gradient estimates for Gradient Ascent

Overview



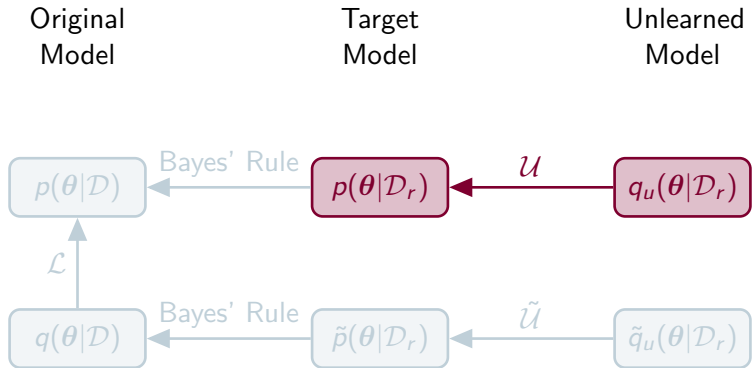
Exact Bayesian Unlearning

- Retraining on \mathcal{D}_r we can get $p(\theta|\mathcal{D}_r)$, but computationally costly,
- Alternatively, using Bayes' rule

$$p(\theta|\mathcal{D}_r) = p(\theta|\mathcal{D}) p(\mathcal{D}_e|\mathcal{D}_r)/p(\mathcal{D}_e|\theta) \propto p(\theta|\mathcal{D})/p(\mathcal{D}_e|\theta)$$

- For discrete θ and conjugate priors can be obtained directly
- Paper focuses on **non-conjugate priors**

Overview



Approximate Bayesian Unlearning

- Find $q_u(\theta|\mathcal{D}_r) \approx p(\theta|\mathcal{D}_r)$ by unlearning from erased data \mathcal{D}_e
- Predictive distributions
 - $q_u(y|\mathcal{D}_r) \triangleq \int p(y|\theta) q_u(\theta|\mathcal{D}_r) d\theta$
 - $p(y|\mathcal{D}_r) = \int p(y|\theta) p(\theta|\mathcal{D}_r) d\theta$
- Loss to minimize: $\text{KL}[q_u(y|\mathcal{D}_r) \parallel p(y|\mathcal{D}_r)]$
 - Closed form may not exist
 - Hard to estimate and optimize

Approximate Bayesian Unlearning

- Find $q_u(\theta|\mathcal{D}_r) \approx p(\theta|\mathcal{D}_r)$ by unlearning from erased data \mathcal{D}_e
- Predictive distributions
 - $q_u(y|\mathcal{D}_r) \triangleq \int p(y|\theta) q_u(\theta|\mathcal{D}_r) d\theta$
 - $p(y|\mathcal{D}_r) = \int p(y|\theta) p(\theta|\mathcal{D}_r) d\theta$
- Loss to minimize: $\text{KL}[q_u(y|\mathcal{D}_r) \parallel p(y|\mathcal{D}_r)]$
 - Closed form may not exist
 - Hard to estimate and optimize

Proposition 1 Bound

$$\text{KL}[q_u(y|\mathcal{D}_r) \parallel p(y|\mathcal{D}_r)] \leq \text{KL}[q_u(\theta|\mathcal{D}_r) \parallel p(\theta|\mathcal{D}_r)]$$

- How to **minimize** $\text{KL}[q_u(\theta|\mathcal{D}_r) \parallel p(\theta|\mathcal{D}_r)]$?

Evidence Upper Bound (EUBO)

- Similar to ELBO define an *evidence upper bound*

$$\mathcal{U} \triangleq \int \underbrace{q_u(\theta|\mathcal{D}_r) \log p(\mathcal{D}_e|\theta) d\theta}_{\text{Forget } \mathcal{D}_e} + \underbrace{\text{KL}[q_u(\theta|\mathcal{D}_r) \parallel p(\theta|\mathcal{D})]}_{\text{Remember } \mathcal{D}(\text{includes } \mathcal{D}_r)}$$

- We can minimize $\text{KL}[q_u(\theta|\mathcal{D}_r) \parallel p(\theta|\mathcal{D}_r)]$ in two ways
 - 1 **Maximize** ELBO while retraining using \mathcal{D}_r
 - 2 **Minimize** EUBO when unlearning with \mathcal{D}_e
- EUBO naturally has regularizing term to avoid *catastrophic forgetting*

EUBO and ELBO

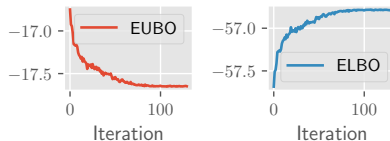
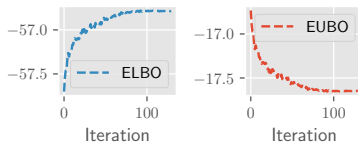
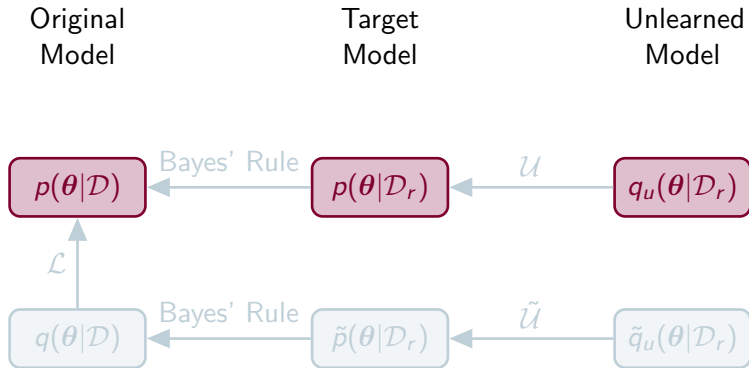
(a) Unlearning from \mathcal{D}_e by minimizing EUBO(b) Retraining with \mathcal{D}_r by maximizing ELBO

Figure: Plots of EUBO and ELBO when (a) unlearning from \mathcal{D}_e and (b) retraining with \mathcal{D}_r .

Overview



Reality Check

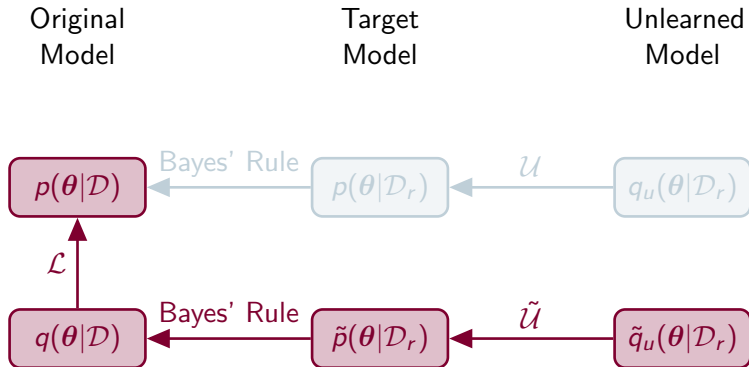
- In reality we only obtain approximations
- VI training gives approximate posterior $q(\boldsymbol{\theta}|\mathcal{D})$
- We can estimate unknown $p(\boldsymbol{\theta}|\mathcal{D}_r)$ using

$$\tilde{p}(\boldsymbol{\theta}|\mathcal{D}_r) \propto q(\boldsymbol{\theta}|\mathcal{D})/p(\mathcal{D}_e|\boldsymbol{\theta})$$

- $\tilde{q}_u(\boldsymbol{\theta}|\mathcal{D}_r)$ from minimizing loss $\text{KL}[\tilde{q}_u(\boldsymbol{\theta}|\mathcal{D}_r) \parallel \tilde{p}(\boldsymbol{\theta}|\mathcal{D}_r)]$
- Define following EUBO

$$\tilde{\mathcal{U}} \triangleq \int \tilde{q}_u(\boldsymbol{\theta}|\mathcal{D}_r) \log p(\mathcal{D}_e|\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta} + \text{KL}[\tilde{q}_u(\boldsymbol{\theta}|\mathcal{D}_r) \parallel q(\boldsymbol{\theta}|\mathcal{D})]$$

Overview



Issues

Two possible sources of inaccuracy in $q(\theta|\mathcal{D})$

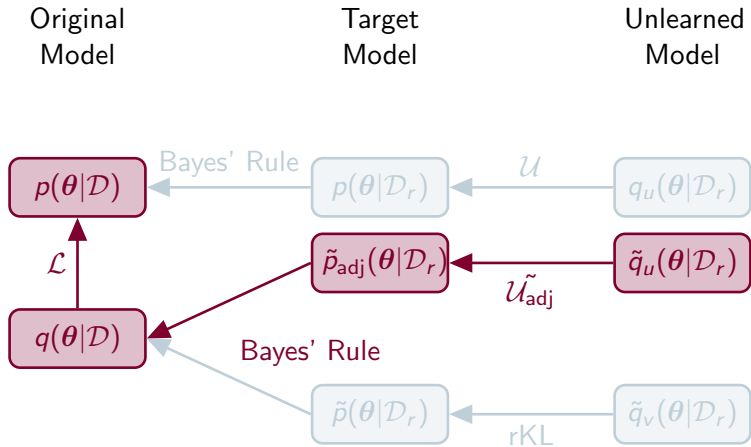
- 1 $q(\theta|\mathcal{D})$ often underestimates the variance of $p(\theta|\mathcal{D})$
- 2 Unlikely that the ELBO is maximized using samples of θ with small $q(\theta|\mathcal{D})$

Hence, curb unlearning at values of θ with small $q(\theta|\mathcal{D})$



Figure: Plot of $q(\theta|\mathcal{D})$ learned using VI. Gray shaded region corresponds to values of θ where $q(\theta|\mathcal{D}) \leq \lambda \max_{\theta'} q(\theta'|\mathcal{D})$. Vertical blue strips on horizontal axis show 100 samples of $\theta \sim q(\theta|\mathcal{D})$.

Overview



Adjusted EUBO

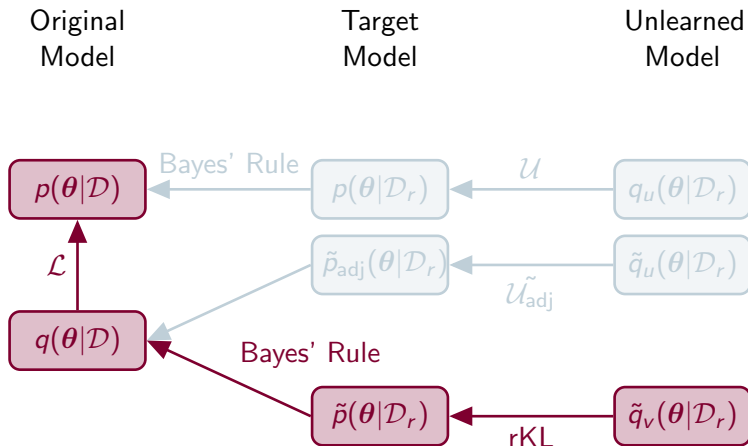
Introduce $\lambda \in [0, 1]$ to control focus of unlearning

$$p_{\text{adj}}(\mathcal{D}_e | \theta; \lambda) \triangleq \begin{cases} p(\mathcal{D}_e | \theta) & \text{if } q(\theta | \mathcal{D}) > \lambda \max_{\theta'} q(\theta' | \mathcal{D}) , \\ 1 & \text{otherwise (i.e., shaded area) ;} \end{cases}$$

$$\tilde{p}_{\text{adj}}(\theta | \mathcal{D}_r; \lambda) \propto \begin{cases} q(\theta | \mathcal{D}) / p(\mathcal{D}_e | \theta) & \text{if } q(\theta | \mathcal{D}) > \lambda \max_{\theta'} q(\theta' | \mathcal{D}) , \\ q(\theta | \mathcal{D}) & \text{otherwise (i.e., shaded area) } \end{cases}$$

$$\tilde{\mathcal{U}}_{\text{adj}}(\lambda) \triangleq \int \tilde{q}_u(\theta | \mathcal{D}_r; \lambda) \log p_{\text{adj}}(\mathcal{D}_e | \theta; \lambda) d\theta + \text{KL}[\tilde{q}_u(\theta | \mathcal{D}_r; \lambda) \parallel q(\theta | \mathcal{D})]$$

Overview



Reverse KL

- Minimize $\text{KL}[\tilde{p}(\boldsymbol{\theta}|\mathcal{D}_r) \parallel \tilde{q}_v(\boldsymbol{\theta}|\mathcal{D}_r)]$ instead
- Now, $\tilde{q}_v(\boldsymbol{\theta}|\mathcal{D}_r)$ **overestimates** the variance of $\tilde{p}(\boldsymbol{\theta}|\mathcal{D}_r)$
- Initialize $\tilde{q}_v(\boldsymbol{\theta}|\mathcal{D}_r)$ at $q(\boldsymbol{\theta}|\mathcal{D})$ for faster convergence
- Now, SGA naturally curbs unlearning at values of $\boldsymbol{\theta}$ with small $q(\boldsymbol{\theta}|\mathcal{D})$

Adjusted EUBO and rKL

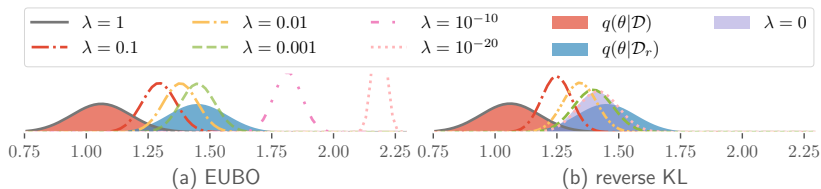
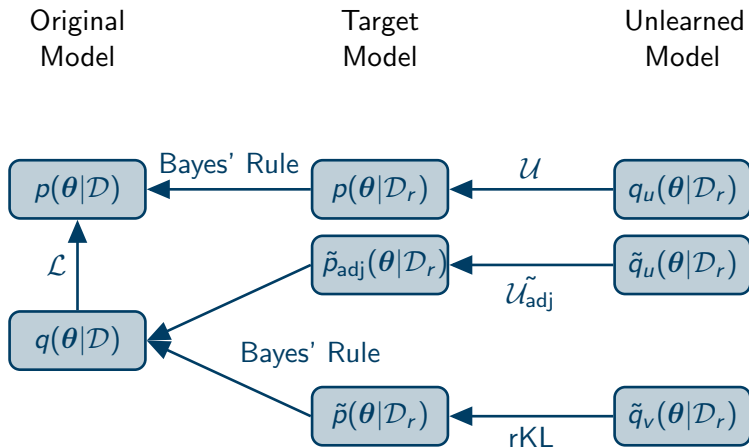


Figure: Plot of approximate posterior beliefs with varying λ obtained by minimizing (a) EUBO (i.e., $\tilde{q}_u(\theta|\mathcal{D}_r; \lambda)$) and (b) reverse KL (i.e., $\tilde{q}_v(\theta|\mathcal{D}_r; \lambda)$); horizontal axis denotes $\theta = \alpha$. In (a), a huge probability mass of $\tilde{q}_u(\theta|\mathcal{D}_r, \lambda = 0)$ is at large values of α beyond the plotting area and the top of the plot of $\tilde{q}_u(\theta|\mathcal{D}_r, \lambda = 10^{-20})$ is cut off due to lack of space.

Overview



Results