

Sample Average Approximation for Black-Box VI

Javier Burroni, Justin Domke, Daniel Sheldon

1. Setup

Focus: Variational inference for statistical models:

- hundreds of variables
- without data-subsampling

ELBO maximization:

- with reparameterizable distribution q_θ via $z_\theta(\cdot)$ and q_{base}

$$\max_{\theta \in \Theta} \mathcal{L}(\theta) = \max_{\theta \in \Theta} \mathbb{E}_{\epsilon \sim q_\theta} \left[\ln \frac{p(z_\theta(\epsilon), x)}{q_\theta(\epsilon)} \right]$$

Usually solved with SGD:

- Hard to tune hyper-parameters
- Results highly dependent on choices

3. SAA

Take $\epsilon_1, \dots, \epsilon_n \sim q_{\text{base}}$

Create deterministic optimization problem:

$$\text{Solve } \max_{\theta \in \Theta} \widehat{\mathcal{L}}_\epsilon(\theta) = \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \left[\ln \frac{p(z_\theta(\epsilon_i), x)}{q_\theta(z_\theta(\epsilon_i))} \right]$$

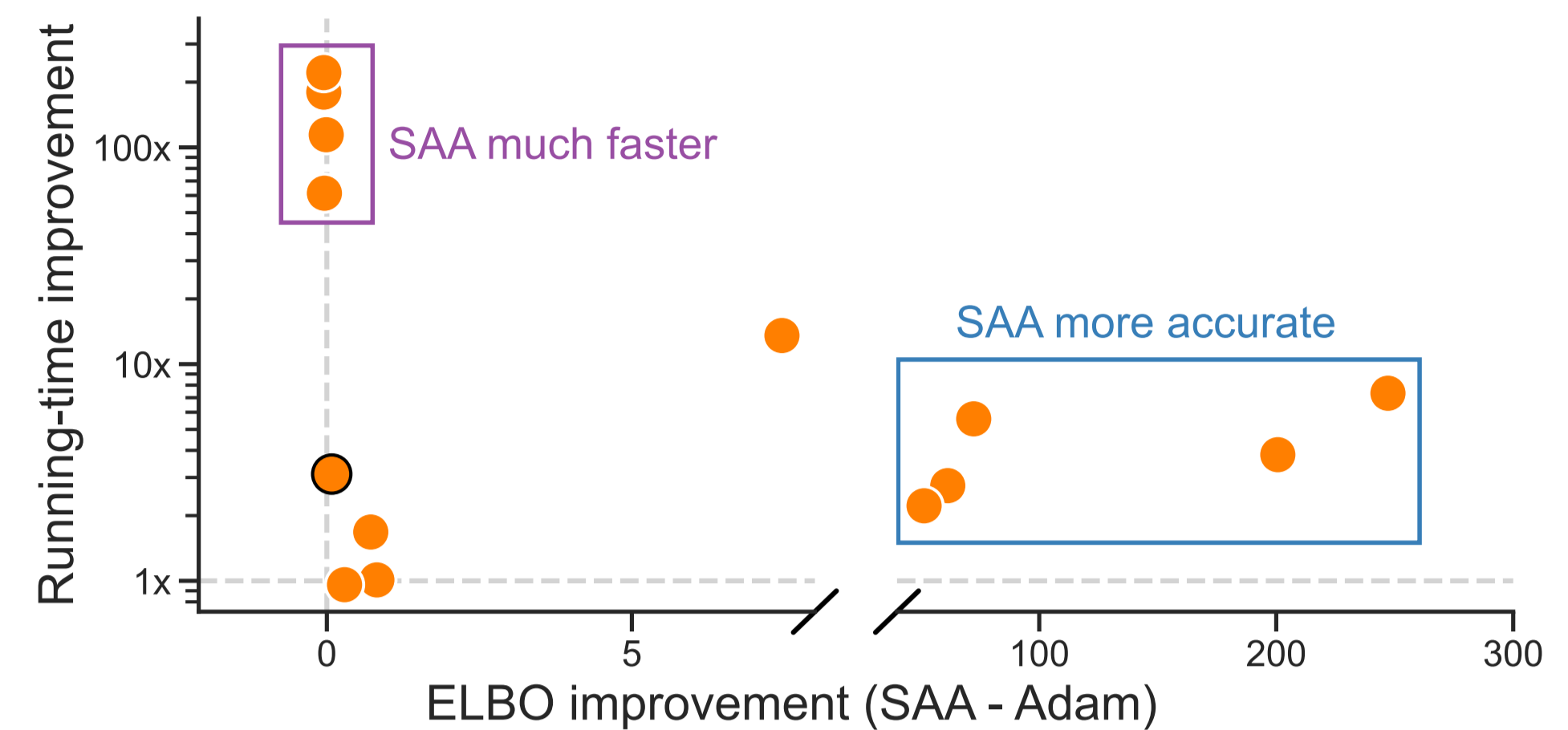
Optimize with:

L-BFGS for search direction and **line search** for step-size

2. Contribution

We introduced **SAA for VI**:

- An alternative stochastic-optimization for BBVI
- Enhances both speed and quality of approximation



4. Sequence of SAA

Use sequence of sizes $n_1 < n_2 < \dots$ to reduce Monte Carlo error

Algorithm 1 SAA for VI

- 1: **while not** converged(θ, ϵ_n) **do**
- 2: $n \leftarrow 2n$
- 3: $\epsilon_n \leftarrow \epsilon_1, \dots, \epsilon_n$, where $\epsilon_i \sim q_{\text{base}}$
- 4: $\theta \leftarrow \text{Optimize}(\theta, \epsilon_n)$
- 5: **return** $\theta^* \leftarrow \theta$

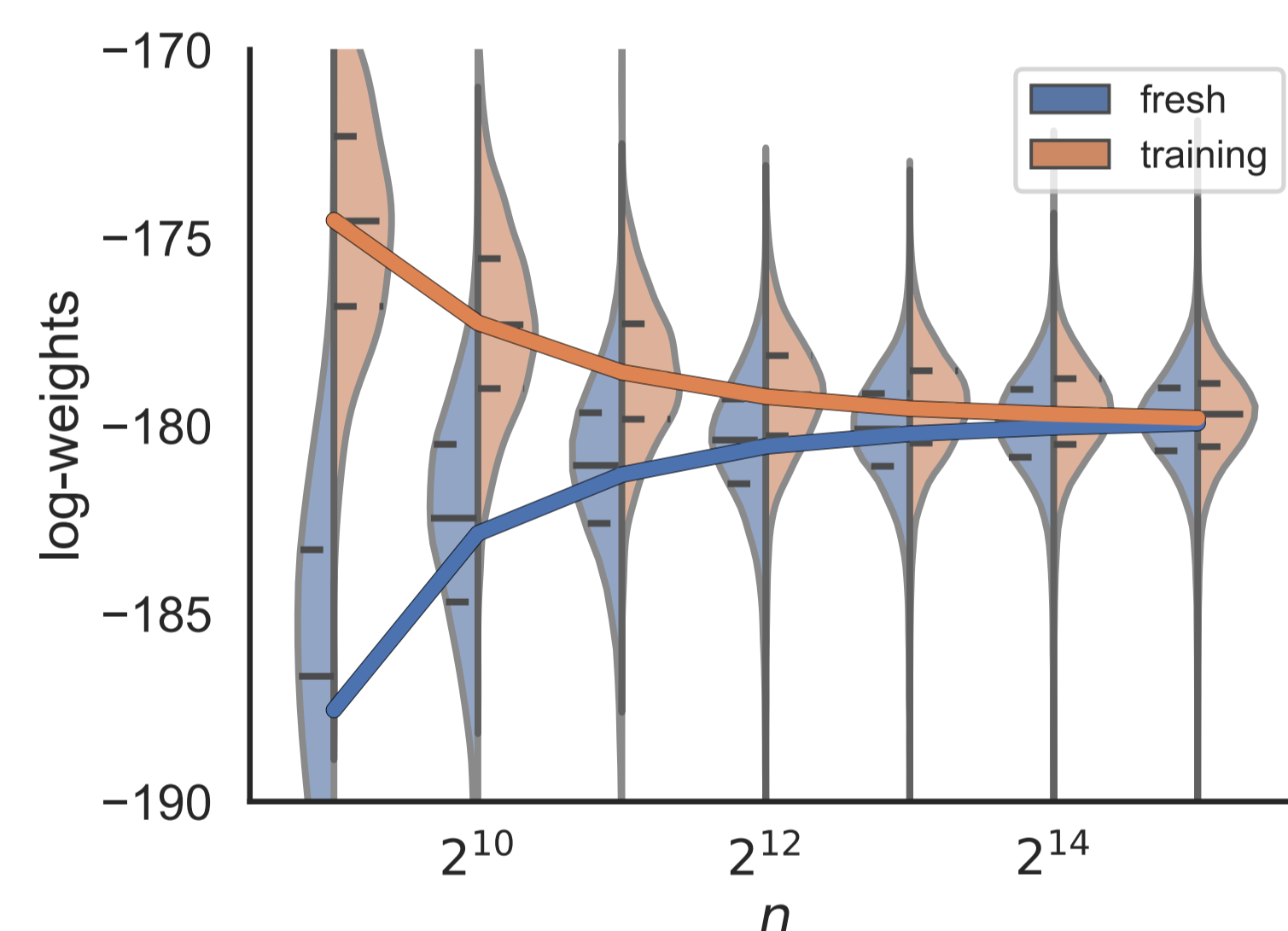
5. New convergence criterion

Compare distributions of log-weights.

Stop when training and testing cannot be distinguished

Training and **testing** log-weights distribution.

By increasing the sample size used for training we achieve a better approximation.

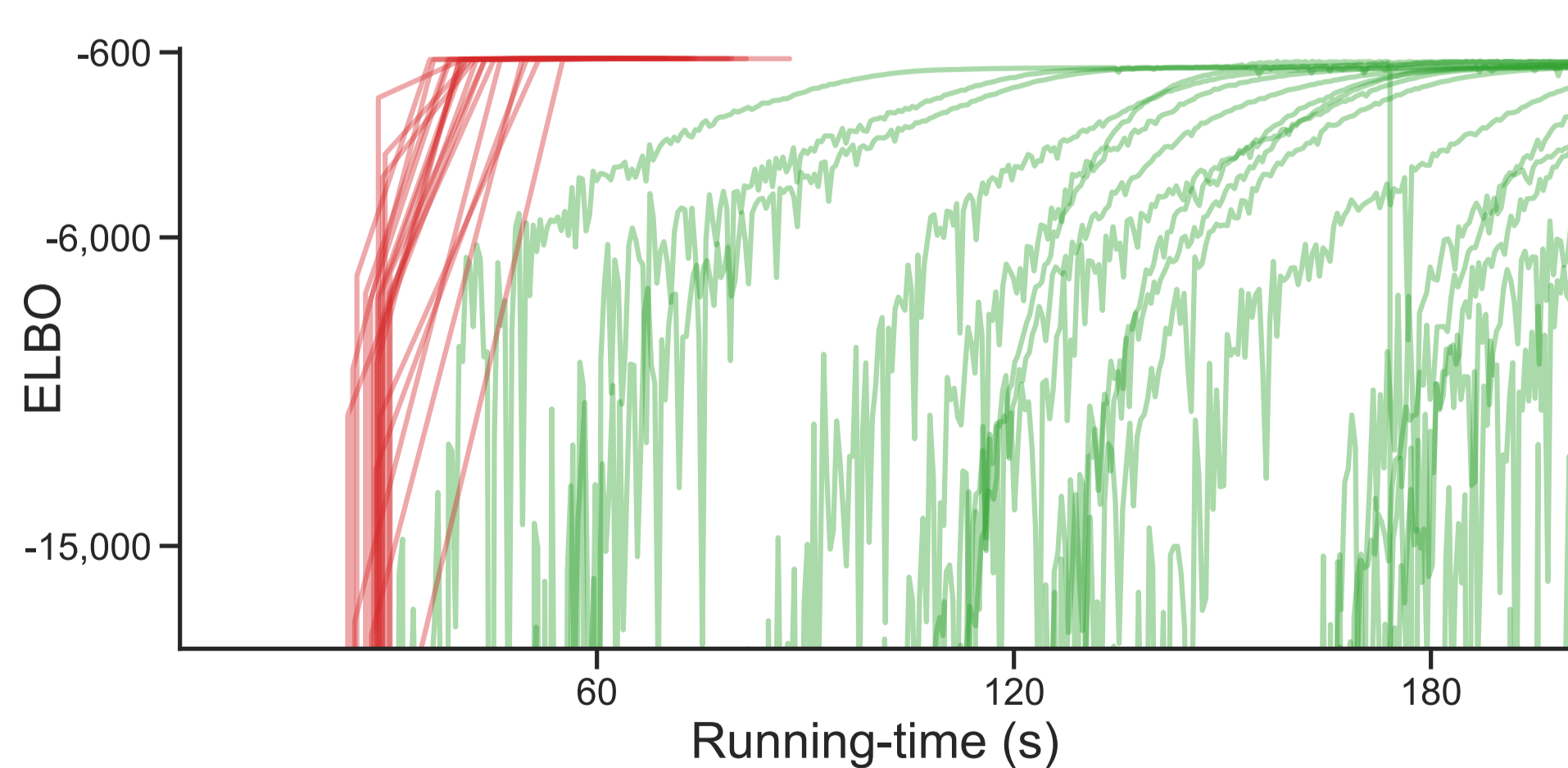
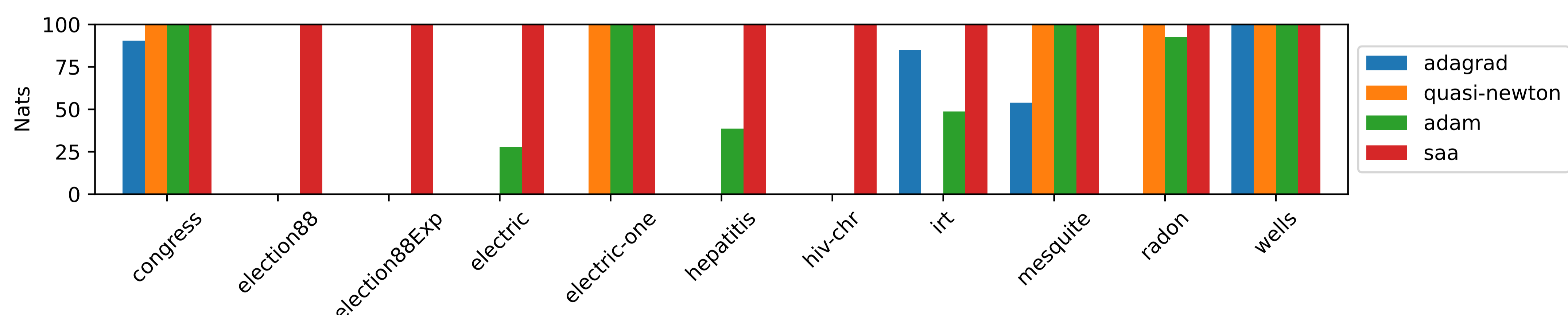


Algorithm 2 converged?

- | | |
|--|---|
| 1: Input: θ, ϵ_n, t | Output: True if converged |
| 2: $\hat{\epsilon}_{10k} \leftarrow \hat{\epsilon}_1, \dots, \hat{\epsilon}_{10k}$, | $\hat{\epsilon}_i \sim q_{\text{base}}$ |
| 3: $\text{obj} \leftarrow \text{mean}(v_\theta(\epsilon_n))$ | |
| 4: $\text{elbo} \leftarrow \text{mean}(v_\theta(\hat{\epsilon}_{10k}))$ | |
| ▷ Statistically compare means: | |
| 5: $p_{\text{value}} \leftarrow t_{\text{test}}(v_\theta(\epsilon_n), v_\theta(\hat{\epsilon}_{10k}))$ | |
| 6: return $p_{\text{value}} > 0.01$ | |

6. Experimental results

Stan models: ELBO comparison after training with dense Gaussian approximation.
For each model, ELBOs are shifted so the best method has value 100.



SAA for VI vs Adam:
Left: electric model ($D = 100$)
Right: Stochastic volatility ($D \approx 17k$)

