

# Kernel-Based Just-In-Time Learning For Passing Expectation Propagation Messages

Wittawat Jitkrittum,<sup>1</sup> Arthur Gretton,<sup>1</sup> Nicolas Heess, S. M. Ali Eslami  
Balaji Lakshminarayanan,<sup>1</sup> Dino Sejdinovic<sup>2</sup> and Zoltán Szabó<sup>1</sup>

Gatsby Unit, University College London<sup>1</sup> University of Oxford<sup>2</sup>

## Introduction

EP is a widely used message passing based inference algorithm.

- **Problem:** Expensive to compute outgoing from incoming messages.
- **Goal:** Speed up computation by a cheap regression function (message operator):

incoming messages  $\mapsto$  outgoing message.

**Merits:**

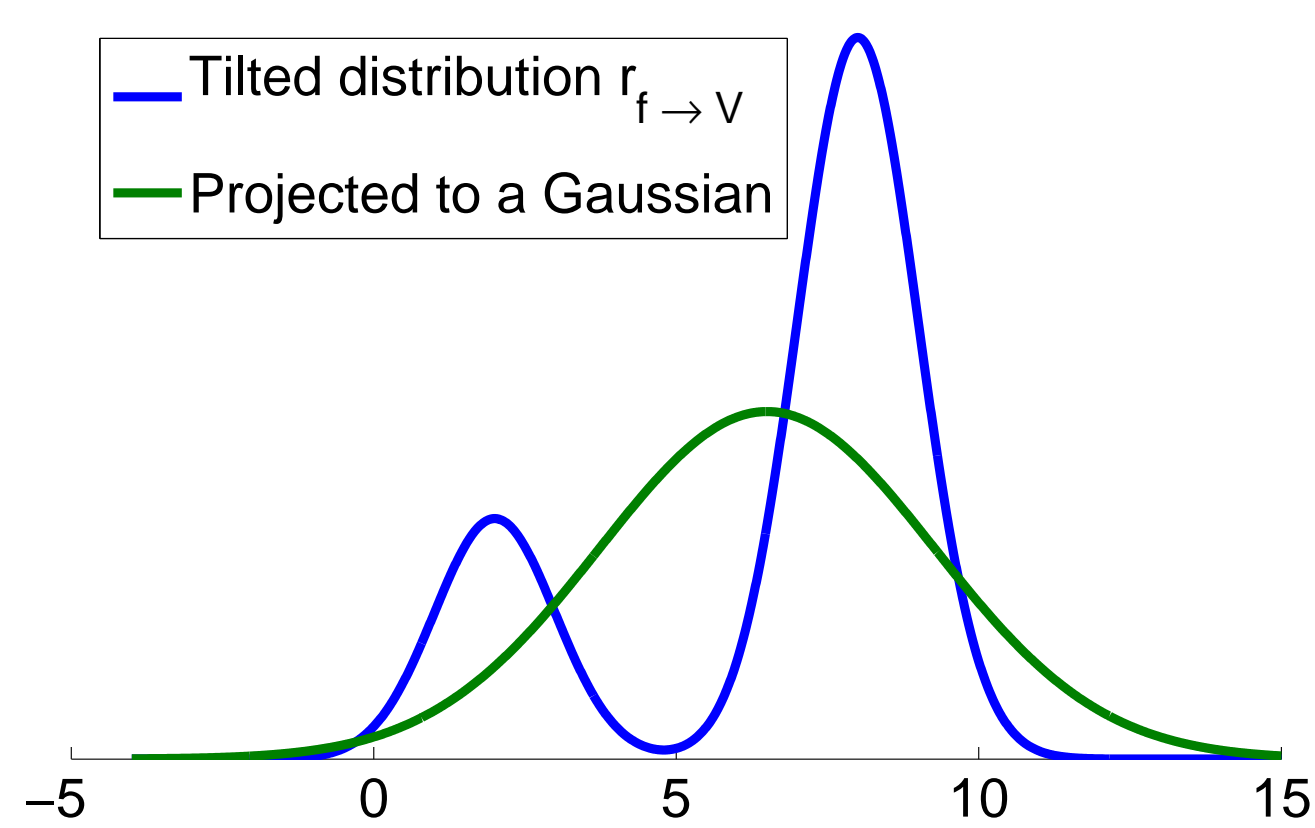
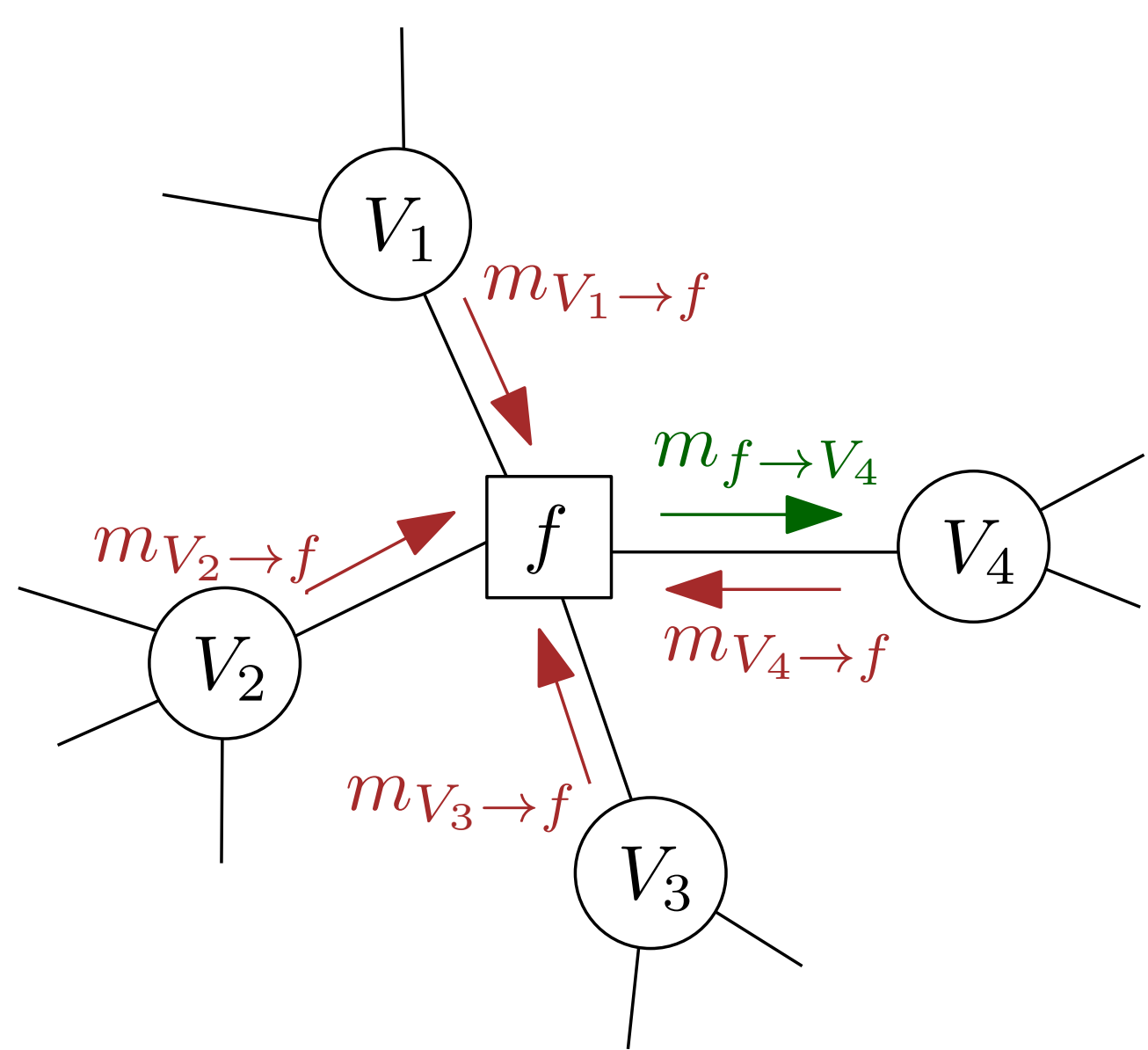
- Efficient online update of the operator during inference.
- Uncertainty monitored to invoke new training examples when needed.
- Automatic random feature representation of incoming messages.

## Expectation Propagation (EP)

Under an approximation that each factor fully factorizes, an outgoing EP message  $m_{f \rightarrow V_i}$  takes the form

$$m_{f \rightarrow V_i}(v_i) = \frac{\text{proj} \left[ f f(\mathcal{V}) \prod_{j=1}^c m_{V_j \rightarrow f}(v_j) d\mathcal{V} \setminus \{v_i\} \right]}{m_{V_i \rightarrow f}(v_i)} := \frac{q_{f \rightarrow V_i}(v_i)}{m_{V_i \rightarrow f}(v_i)}$$

$\text{proj}[r_{f \rightarrow V_i}] := \arg \min_{q \in \text{ExpFam}} \text{KL}[r_{f \rightarrow V_i} \parallel q]$   
(projection onto exponential family)



**Projected message:**

- $q_{f \rightarrow V}(v) = \text{proj}[r_{f \rightarrow V}(v)] \in \text{ExpFam}$  with sufficient statistic  $u(v)$ .
- Compute  $q_{f \rightarrow V}(v)$  by moment matching:  $\mathbb{E}_{q_{f \rightarrow V}}[u(v)] = \mathbb{E}_{r_{f \rightarrow V}}[u(v)]$ .

## Kernel on Incoming Messages

Propose to incrementally learn during inference a kernel-based EP message operator (distribution-to-distribution regression)

$$[m_{V_j \rightarrow f}]_{j=1}^c \mapsto q_{f \rightarrow V_i}$$

for any factor  $f$  that can be sampled.

- Product distribution of  $c$  incoming messages:  $\mathbf{r} := \times_{l=1}^c r_l$ ,  $\mathbf{s} := \times_{l=1}^c s_l$ .
- Mean embedding of  $\mathbf{r}$ :  $\mu_r := \mathbb{E}_{a \sim r} k(\cdot, a)$ .
- Gaussian kernel on (product) distributions. Two-staged random feature approx.:

$$\kappa(\mathbf{r}, \mathbf{s}) = \exp\left(-\frac{\|\mu_r - \mu_s\|_{\mathcal{H}}^2}{2\gamma^2}\right) \approx \exp\left(-\frac{\|\hat{\phi}(\mathbf{r}) - \hat{\phi}(\mathbf{s})\|_{D_{\text{in}}}^2}{2\gamma^2}\right) \approx \hat{\psi}(\mathbf{r})^\top \hat{\psi}(\mathbf{s}).$$

## Message Operator: Bayesian Linear Regression

- **Input:**  $\mathbf{X} = (\mathbf{x}_1 | \dots | \mathbf{x}_N)$ :  $N$  training incoming messages represented as random feature vectors.
- **Output:**  $\mathbf{Y} = (\mathbb{E}_{r_1 \rightarrow V} u(v) | \dots | \mathbb{E}_{r_N \rightarrow V} u(v)) \in \mathbb{R}^{D_y \times N}$ : expected sufficient statistics of outgoing messages.
- Inexpensive online updates of posterior mean and covariance.
- Bayesian regression gives prediction and predictive variance.
- If predictive variance  $>$  threshold, query the importance sampling oracle.

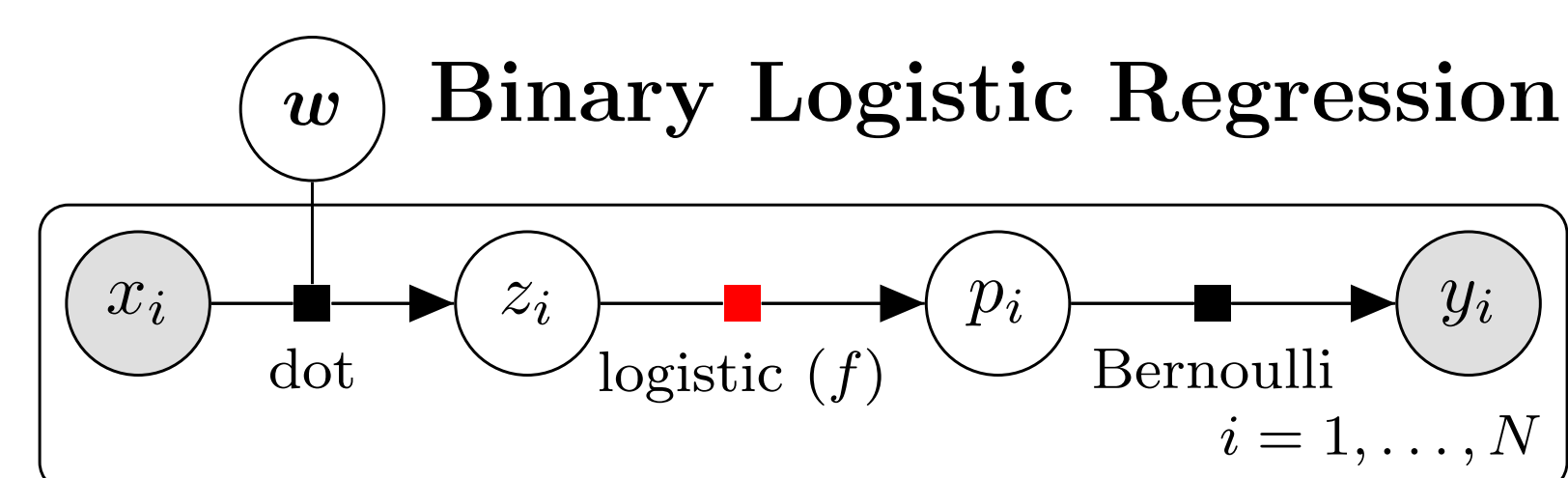
## Two-Staged Random Features

**In:**  $\mathcal{F}(k)$ : Fourier transform of  $k$ ,  $D_{\text{in}}$ : #inner features,  $D_{\text{out}}$ : #outer features,  $k_{\text{gauss}}$ : Gaussian kernel on  $\mathbb{R}^{D_{\text{in}}}$

**Out:** Random features  $\hat{\psi}(\mathbf{r}) \in \mathbb{R}^{D_{\text{out}}}$

- 1: Sample  $\{\omega_i\}_{i=1}^{D_{\text{in}}} \stackrel{i.i.d.}{\sim} \mathcal{F}(k)$ ,  $\{b_i\}_{i=1}^{D_{\text{in}}} \stackrel{i.i.d.}{\sim} U[0, 2\pi]$ .
- 2:  $\hat{\phi}(\mathbf{r}) = \sqrt{\frac{2}{D_{\text{in}}}} (\mathbb{E}_{x \sim r} \cos(\omega_i^\top x + b_i))_{i=1}^{D_{\text{in}}} \in \mathbb{R}^{D_{\text{in}}}$
- 3: Sample  $\{\nu_i\}_{i=1}^{D_{\text{out}}} \stackrel{i.i.d.}{\sim} \mathcal{F}(k_{\text{gauss}}(\gamma^2))$ ,  $\{c_i\}_{i=1}^{D_{\text{out}}} \stackrel{i.i.d.}{\sim} U[0, 2\pi]$ .
- 4:  $\hat{\psi}(\mathbf{r}) = \sqrt{\frac{2}{D_{\text{out}}}} (\cos(\nu_i^\top \hat{\phi}(\mathbf{r}) + c_i))_{i=1}^{D_{\text{out}}} \in \mathbb{R}^{D_{\text{out}}}$

## Experiment 1: Uncertainty Estimates

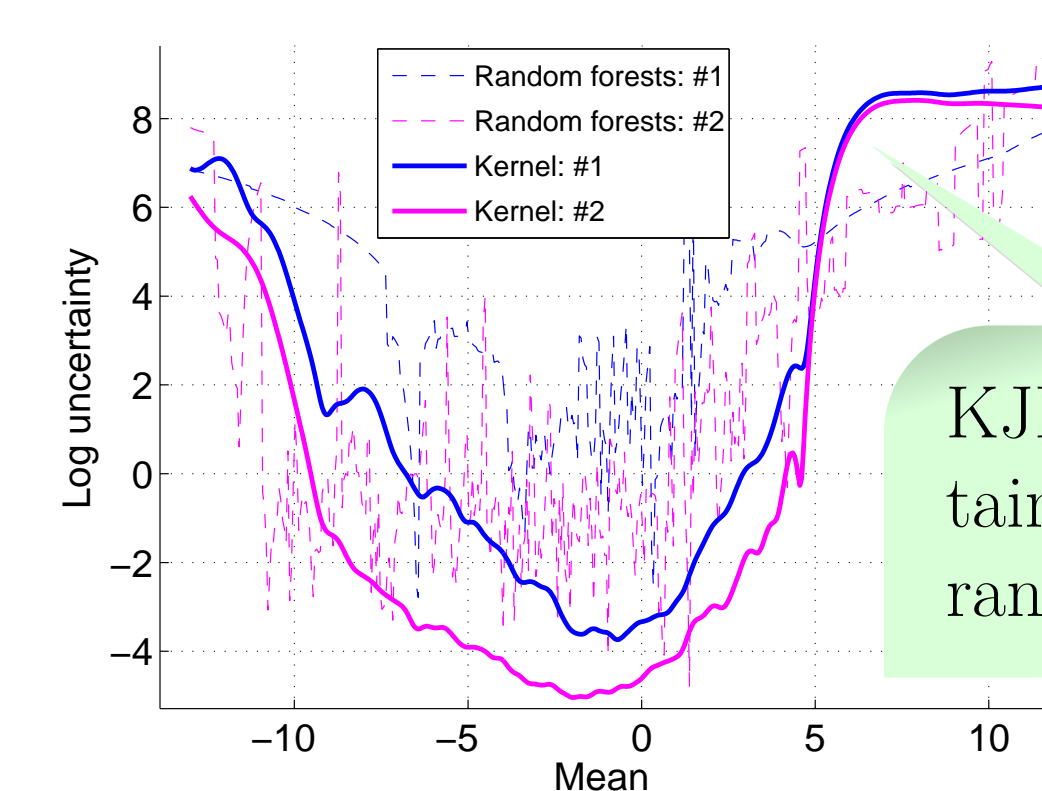
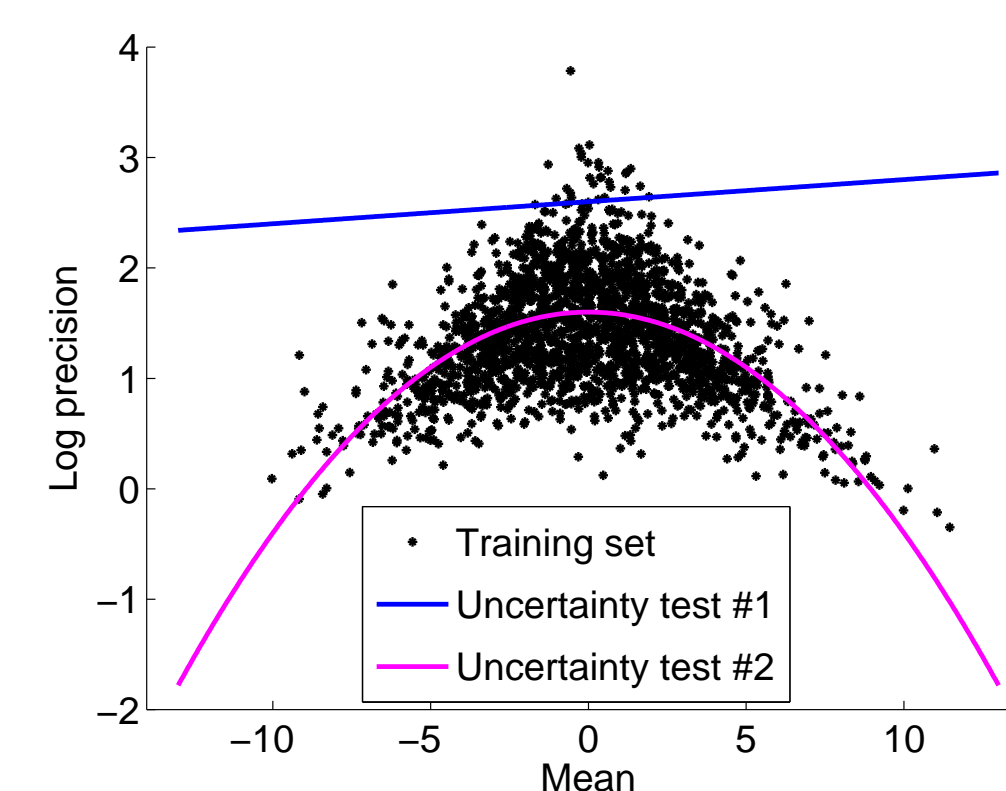


- Approx.  $f(p|z) = \delta\left(p - \frac{1}{1+\exp(-z)}\right)$ .
- Incoming messages:

$$m_{z_i \rightarrow f} = \mathcal{N}(z_i; \mu, \sigma^2),$$

$$m_{p_i \rightarrow f} = \text{Beta}(p_i; \alpha, \beta).$$

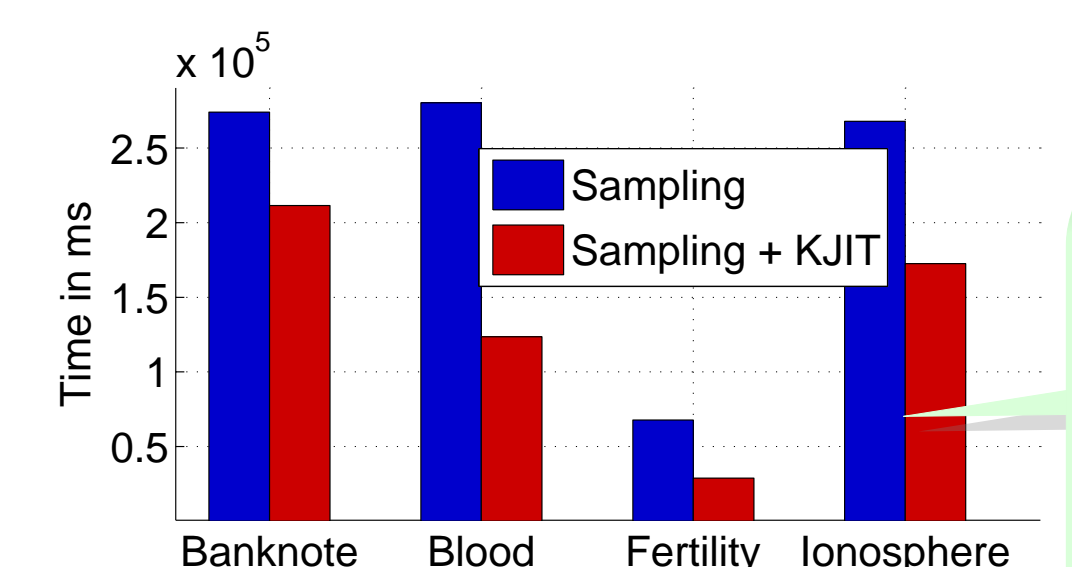
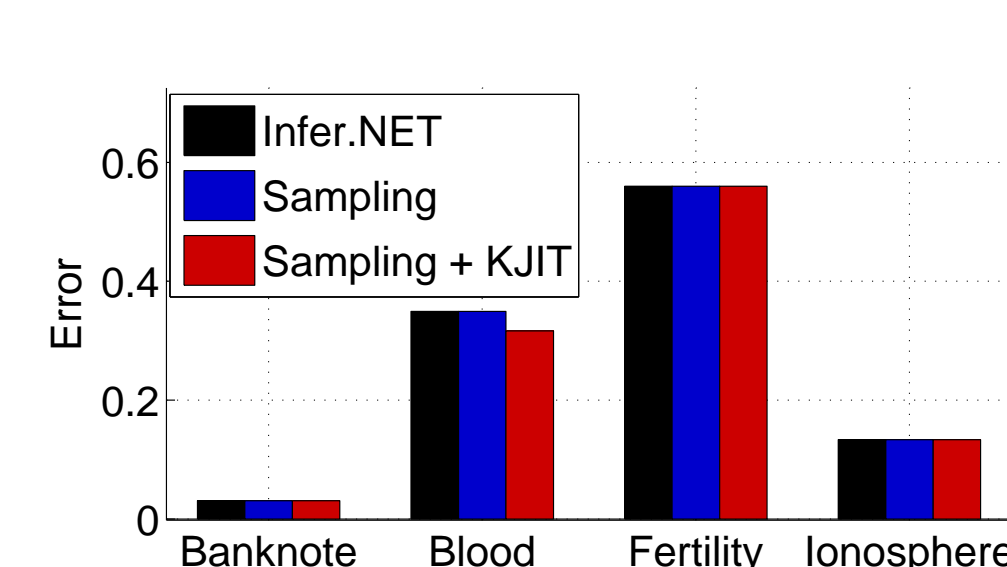
- Training messages collected from 20 EP runs on toy data.
- #Random features:  $D_{\text{in}} = 300$  and  $D_{\text{out}} = 500$ .



KJIT gives smoother uncertainty estimates compared to random forests.

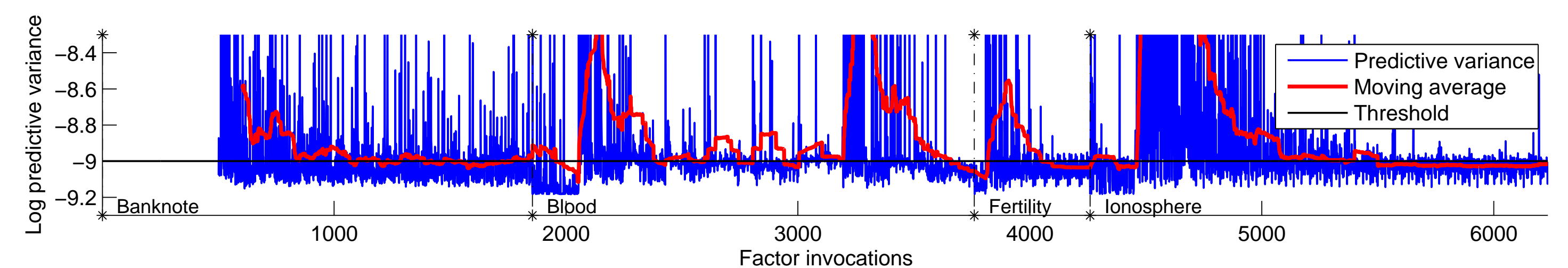
## Experiment 2: Real Data

- Binary logistic regression. Sequentially present 4 real datasets to the operator.
- Diverse distributions of incoming messages.



Much faster with same classification errors as obtained by importance sampling.

- **Sampling + KJIT** = proposed KJIT with an importance sampling oracle.
- KJIT operator can adapt to the change of input message distributions.



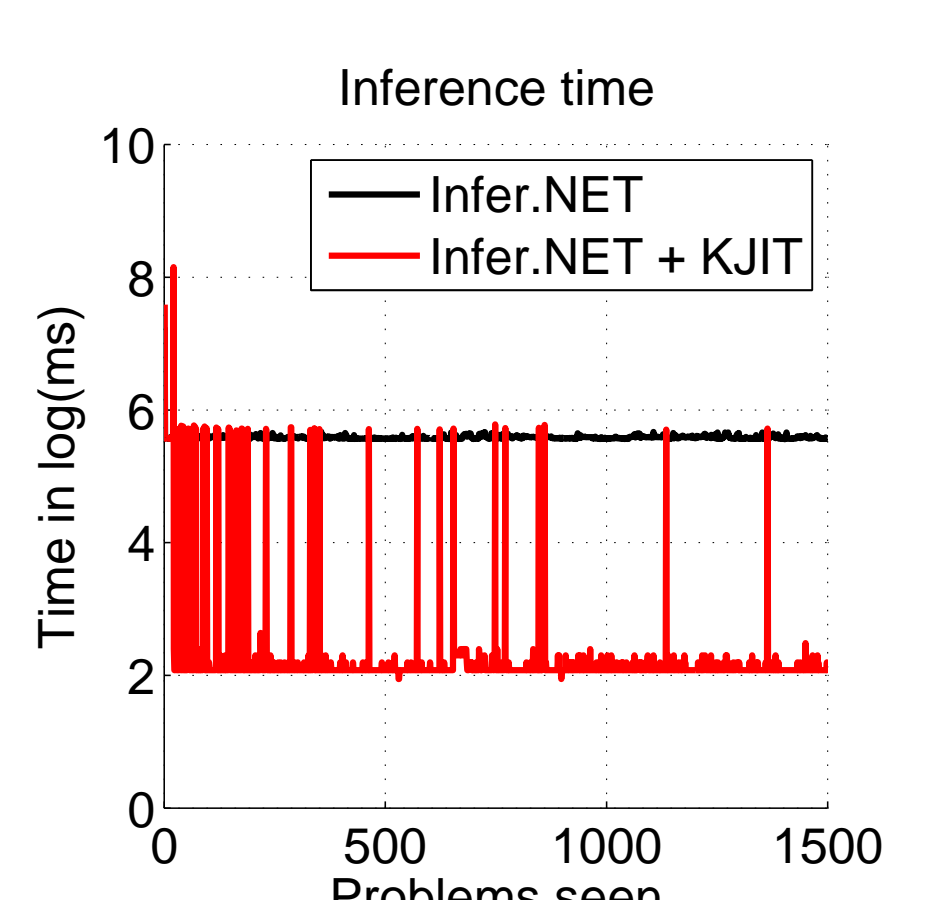
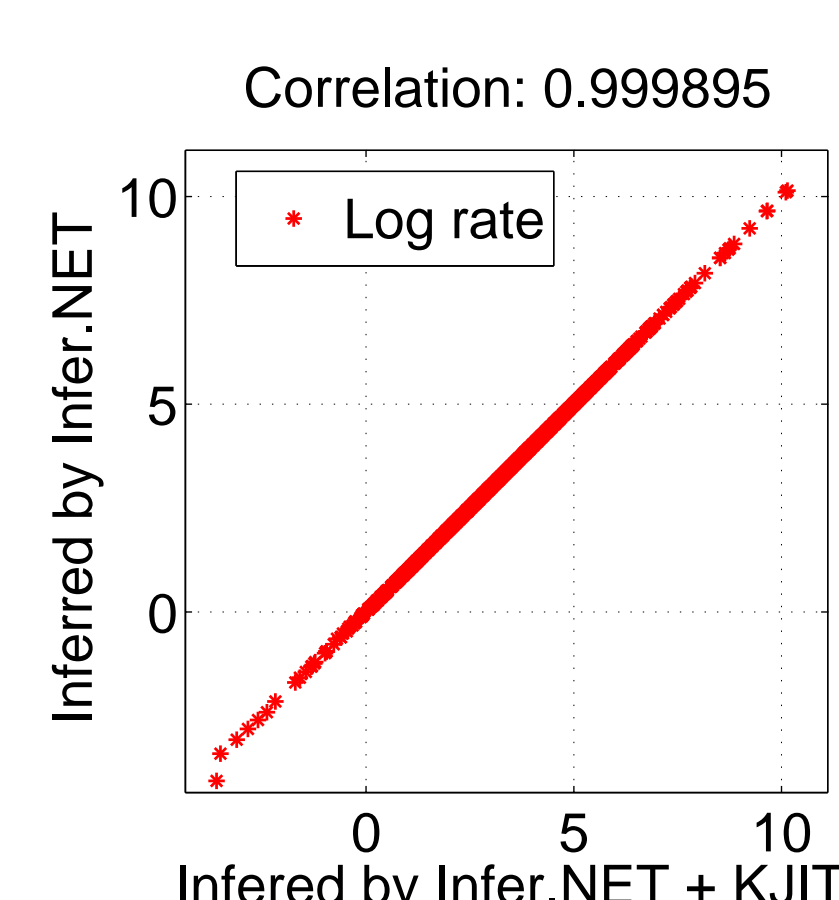
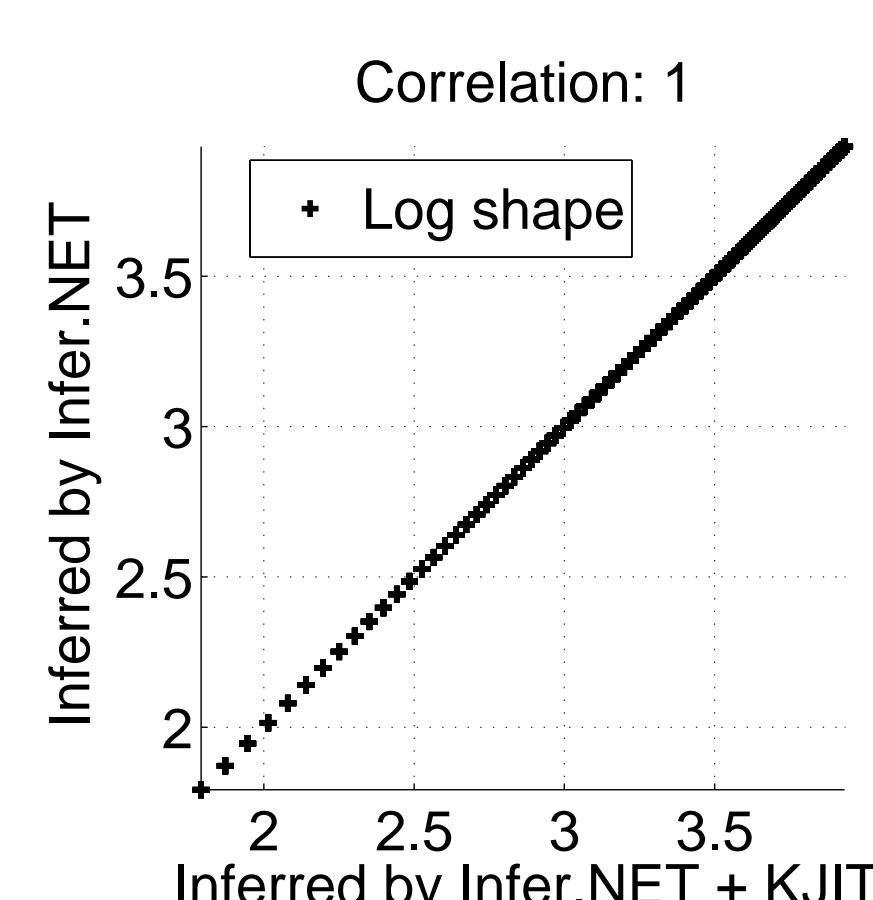
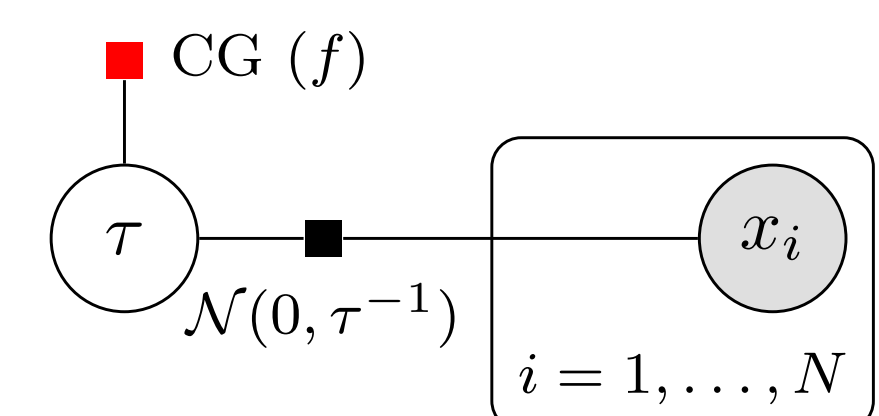
## Experiment 3: Compound Gamma Factor

Infer posterior of the precision  $\tau$  of  $x \sim \mathcal{N}(x; 0, \tau^{-1})$  from observations  $\{x_i\}_{i=1}^N$ :

$$r_2 \sim \text{Gamma}(r_2; s_1, r_1)$$

$$\tau \sim \text{Gamma}(\tau; s_2, r_2)$$

$$(s_1, r_1, s_2) = (1, 1, 1).$$



- **Infer.NET + KJIT** = proposed KJIT with a hand-crafted factor as oracle.
- **Inference quality:** as good as hand-crafted factor; much faster.