

Kernel-Based Just-In-Time Learning for Passing Expectation Propagation Messages

Wittawat Jitkrittum¹, Arthur Gretton¹, Nicolas Heess, S. M. Ali Eslami, Balaji Lakshminarayanan¹, Dino Sejdinovic², and Zoltán Szabó¹
 Gatsby Computational Neuroscience Unit, University College London¹ University of Oxford²



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):

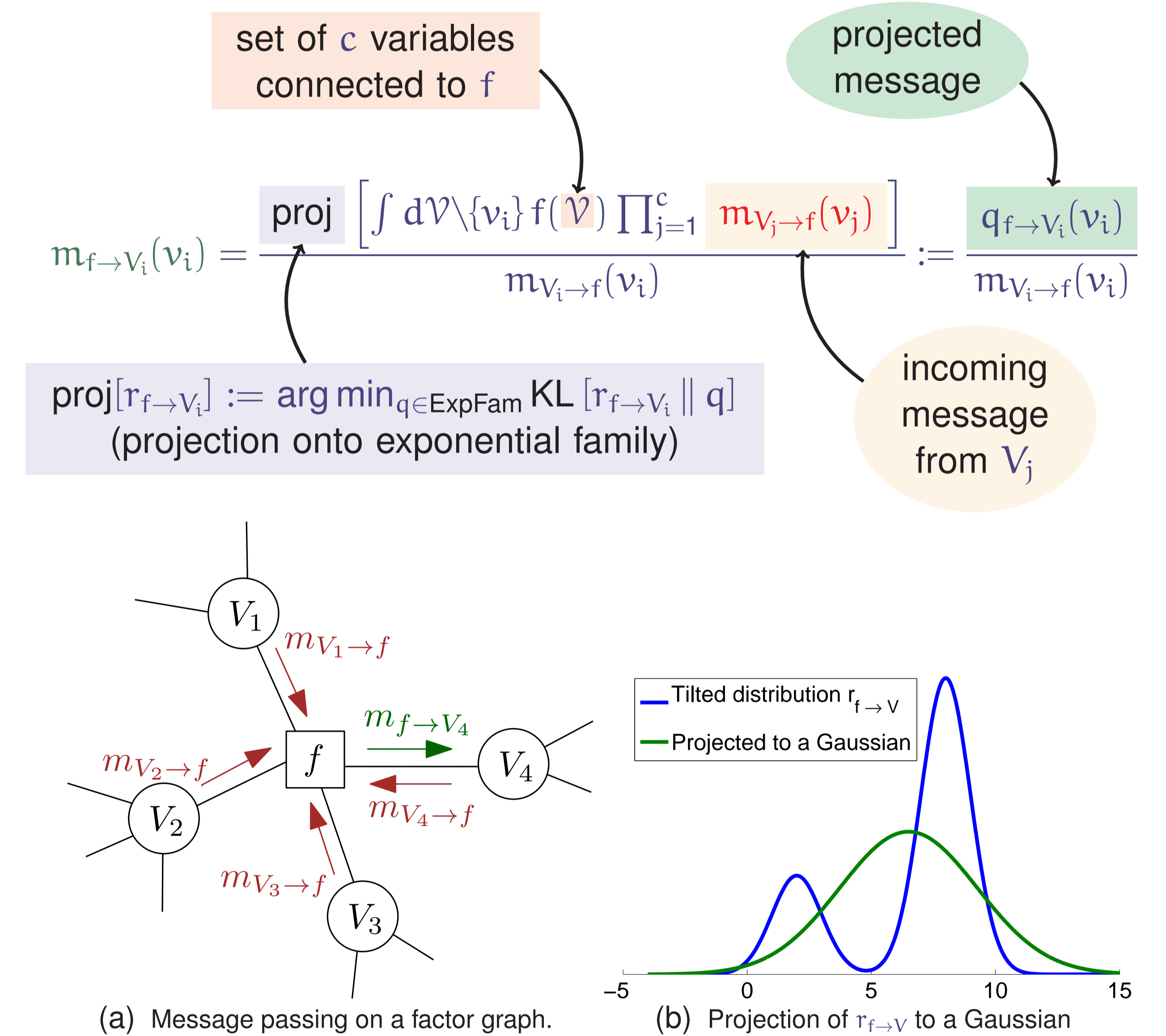
incoming messages \mapsto outgoing message.

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



Projected message:

- $q_{f \rightarrow V}(v) = \text{proj} [r_{f \rightarrow V}(v)] \in \text{ExpFam}$ with sufficient statistic $u(v)$.
- 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: $r := \times_{l=1}^c r_l$, $s := \times_{l=1}^c s_l$.
- Mean embedding of r : $\mu_r := \mathbb{E}_{a \sim r}(\cdot, a)$.
- Gaussian kernel on (product) distributions:

$$\kappa(r, s) = \exp \left(-\frac{\|\mu_r - \mu_s\|_{\mathcal{H}}^2}{2\gamma^2} \right).$$

Two-staged random feature approximation:

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

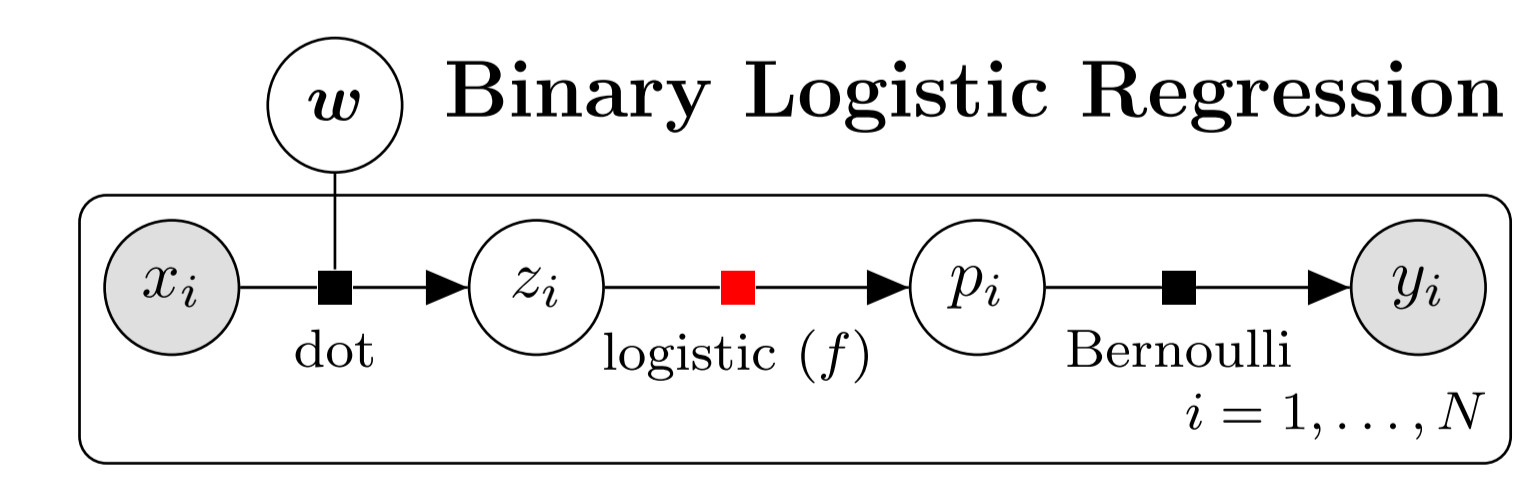
Message Operator: Bayesian Linear Regression

- Input:** $X = (x_1 | \dots | x_N)$: N training incoming messages represented as random feature vectors.
- Output:** $Y = (\mathbb{E}_{r_1 \rightarrow V} u(v) | \dots | \mathbb{E}_{r_N \rightarrow V} u(v)) \in \mathbb{R}^{D_y \times N}$: sufficient statistics of outgoing messages.
- Inexpensive online update.
- Bayesian regression gives prediction and predictive variance.
- If predictive variance $<$ threshold, query importance sampling oracle.

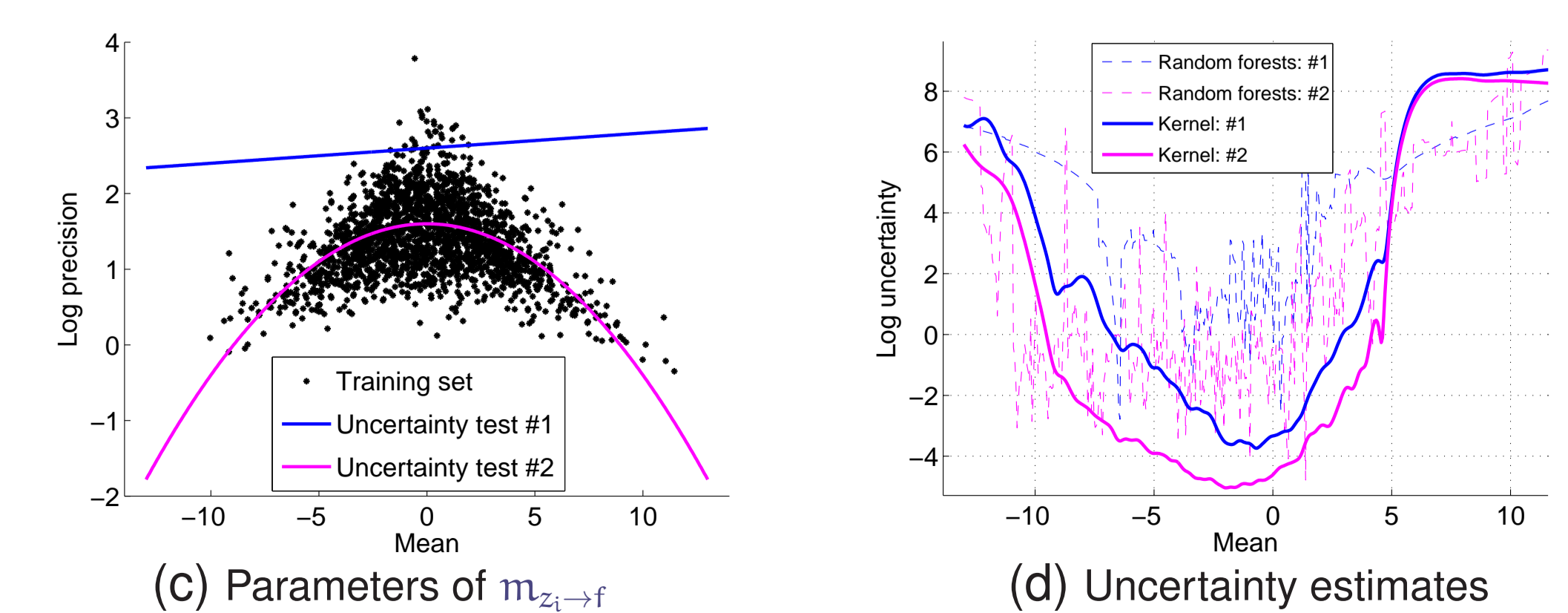
Two-Stage Random Features

- In:** $\mathcal{F}(k)$: Fourier transform of k , D_{in} : #inner features, D_{out} : #outer features, k_{gauss} : Gaussian kernel on $\mathbb{R}^{D_{in}}$
- Out:** Random features $\hat{\psi}(r) \in \mathbb{R}^{D_{out}}$
- Sample $\{\omega_{ij}\}_{i=1}^{D_{in}} \stackrel{i.i.d.}{\sim} \mathcal{F}(k)$, $\{b_{ij}\}_{i=1}^{D_{in}} \stackrel{i.i.d.}{\sim} U[0, 2\pi]$.
 - $\hat{\phi}(r) = \sqrt{\frac{2}{D_{in}}} (\mathbb{E}_{x \sim r} \cos(\omega_i^\top x + b_i))_{i=1}^{D_{in}} \in \mathbb{R}^{D_{in}}$
 - Sample $\{v_{ij}\}_{i=1}^{D_{out}} \stackrel{i.i.d.}{\sim} \mathcal{F}(k_{\text{gauss}}(\gamma^2))$, $\{c_{ij}\}_{i=1}^{D_{out}} \stackrel{i.i.d.}{\sim} U[0, 2\pi]$.
 - $\hat{\psi}(r) = \sqrt{\frac{2}{D_{out}}} (\cos(v_i^\top \hat{\phi}(r) + c_i))_{i=1}^{D_{out}} \in \mathbb{R}^{D_{out}}$

Experiment 1: Uncertainty Estimates



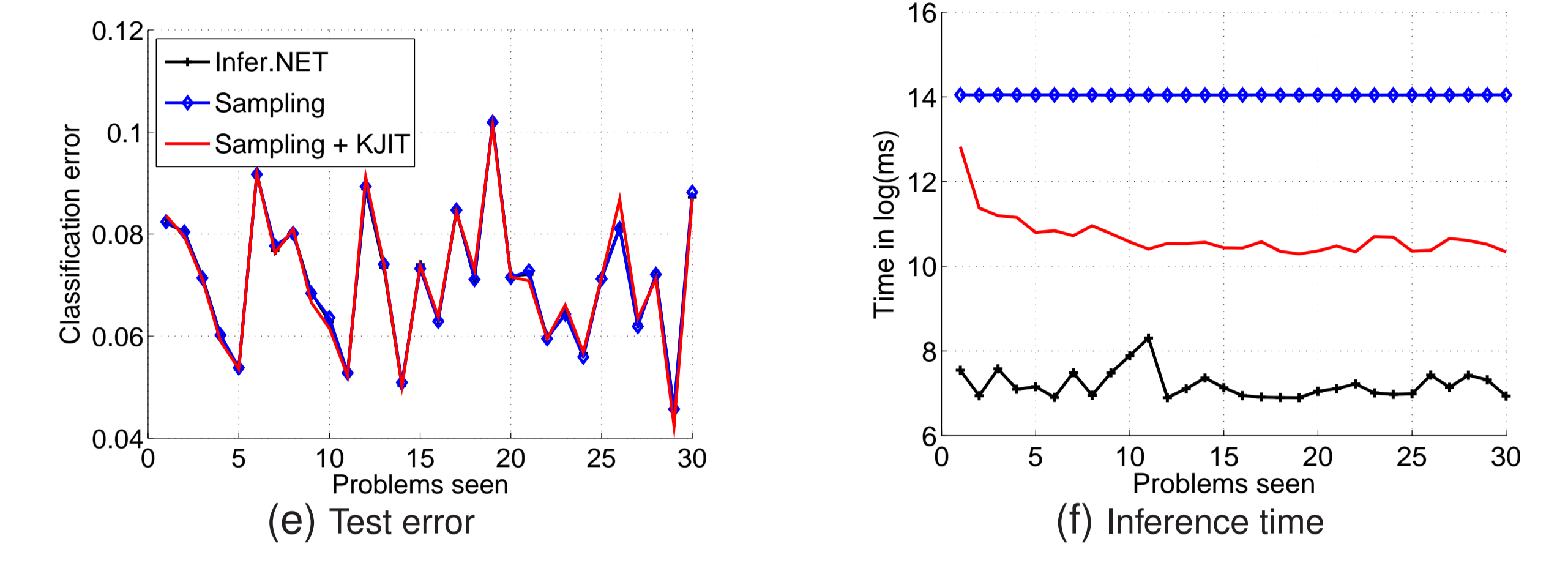
- Approximate the logistic factor: $f(z|x) = \delta \left(z - \frac{1}{1+\exp(-x)} \right)$.
- Incoming messages: $m_{z_i \rightarrow f} = \mathcal{N}(z_i; \mu, \sigma^2)$, $m_{p_i \rightarrow f} = \text{Beta}(p_i; \alpha, \beta)$.
- Training set = messages collected from 20 EP runs on toy data.



#Random features: $D_{in} = 300$ and $D_{out} = 500$.

Experiment 2: Classification Errors

Fix true w . Sequentially present 30 problems. Generate $\{(x_i, y_i)\}_{i=1}^{300}$ for each.



Sampling + KJIT = proposed KJIT with an importance sampling oracle.

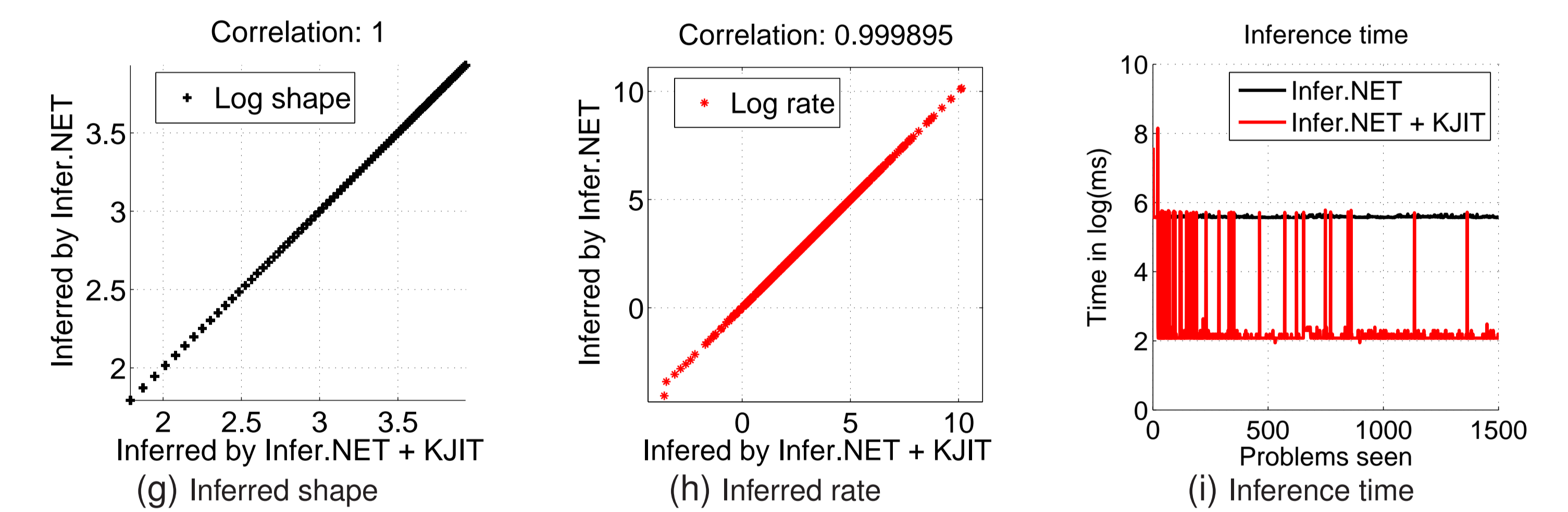
Experiment 3: Compound Gamma Factor

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

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

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

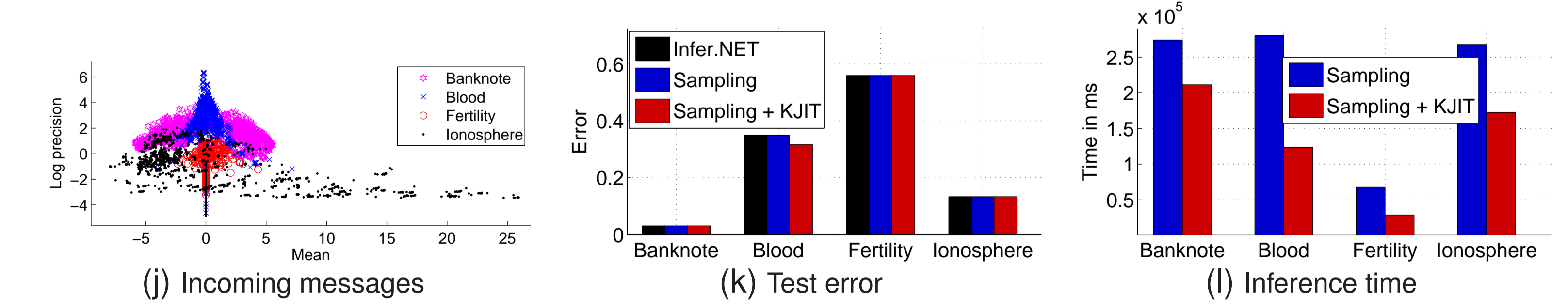
$$(s_1, r_1, s_2) = (1, 1, 1).$$



Inference quality: as good as hand-crafted factor; much faster.

Experiment 4: Real Data

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



KJIT operator can adapt to the change of input message distributions.

