



Neural Network Online Training with Sensitivity to Multiscale Temporal Structure

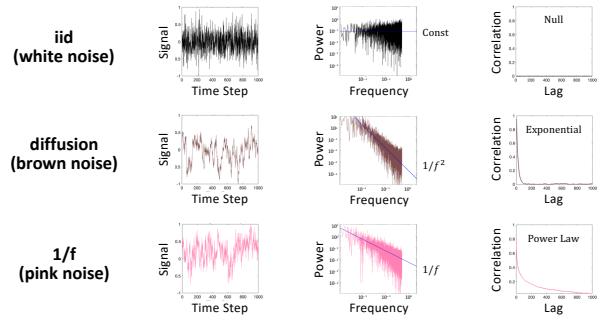
Matt Jones, David Mayo, Tyler R. Scott, Mengye Ren, Gamaleldin ElSayed, Katherine Hermann, Michael C. Mozer
mcjones@google.com

Contributions

- Realistic nonstationarity as $1/f$ noise in data-generating parameters
- Kalman filter for Bayesian inference over $1/f$ noise
- Extended Kalman filter for NN optimization
- Efficient variational approximation, amounts to subweights with different learning and decay rates
- Variational multiscale EKF learns well in nonstationary tasks

Temporal structure in realistic environments

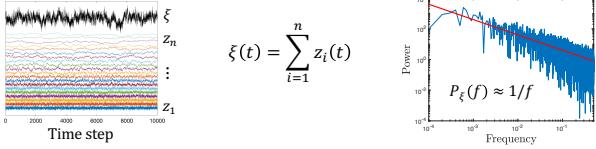
- Real environments have temporal structure (not iid)
- Have long-range autocorrelations (not diffusion or random walk)



Bayesian inference over $1/f$ noise

- Generative model: sum of diffusion processes at different timescales
- Inferential model: Kalman filter over all processes jointly [KTS07]

Generative Model



Timescales $\tau_i = v^i$ ($v > 1$)

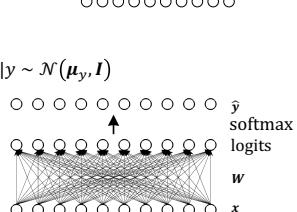
Inferential Model

$$\begin{aligned} Z(t)|\xi_{\leq t} &\sim \mathcal{N}(\omega(t), S(t)) \\ \omega(t+1) &= D \left(\omega(t) - \frac{S(t)}{1^T S(t) 1} \mathbf{1} (\hat{\xi}(t) - \xi(t)) \right) \\ \text{Decay: } D &= \text{diag}(e^{-1/\tau}) \quad \text{Preconditioner} \end{aligned}$$

Loss gradient: $\mathcal{L} = \frac{1}{2} (\hat{\xi} - \xi)^2$

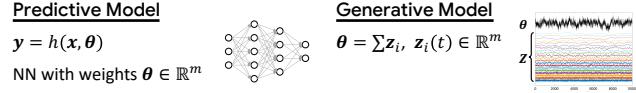
Evaluation in synthetic nonstationary tasks

- Linear regression
 - Data sampling: $y = x^\top \beta$, $\beta \sim 1/f$ noise
 - Predictive model: $y = x^\top \theta$
- 10-way classification
 - Data sampling: $y \sim \text{softmax}(l)$, $x|y \sim \mathcal{N}(\mu_y, I)$
 - Predictive model: $\hat{y} = \text{softmax}(x^\top W)$



Extended Kalman filter and variational inference

- Assume $1/f$ noise in optimal network parameters
- Extended Kalman filter uses Jacobian to linearize network
- Variational approximation tracks variance of each component

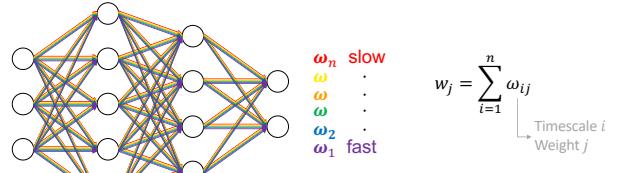


Inferential Model $Z = (z_1, \dots, z_n)$ $Z x_{<t}, y_{<t} \sim \mathcal{N}(\omega(t), S(t))$	Extended Kalman filter (EKF) [SW89, PF03] $h(x(t), \theta) \approx \hat{y} + J_h(\omega(t))(\theta - \omega(t))$ $y \hat{y} \sim \mathcal{N}(\hat{y}, R_{\hat{y}}), R_{\hat{y}} = \text{Var}(y \hat{y})$
---	---

Variational approximation $S(t) \approx \text{diag}(s^2(t))$ Linear in NN size	
---	--

Variational EKF optimizer

- Each weight decomposed as sum of subweights
- Subweight ω_{ij} is mean estimate of θ_j at timescale i (z_{ij})
- Decay rates come from $1/f$ generative model
- Learning rates come from EKF posterior variance



Mean update

$$\omega_{ij}(t+1) = \underbrace{e^{-1/\tau_i}}_{\text{Decay}} \left(\omega_{ij}(t) - \underbrace{\bar{s}_{ij}^2(t)}_{\text{Fisher precision information}} \underbrace{\partial_{w_j} \bar{\mathcal{L}}(t)}_{\text{Loss gradient}} \right)$$

Variance update

$$\bar{s}_{ij}^{-2}(t) = \underbrace{J^\top R^{-1} J}_{\text{Posterior}}_{ij,ij} + \underbrace{e^{-2/\tau_i} \bar{s}_{ij}^{-2}(t-1)}_{\text{Previous Decay posterior}} + \underbrace{2\rho^2(1-e^{-2/\tau_i})}_{\text{Prior precision}}$$

Training methods

- Window: use past H observations
- SGD: update after each observation
- Discounting: power-law weighting
- Variational EKF
- Full EKF

