

# Extended Kalman Filter (EKF) for MLP Training: Derivation and Implementation

## Standard Kalman Filter

The **Kalman filter** solves the linear Gaussian state-estimation problem. With a linear state transition and observation model,

$$\begin{aligned}x_k &= A x_{k-1} + w_{k-1}, \quad w \sim \mathcal{N}(0, Q), \\y_k &= H x_k + v_k, \quad v \sim \mathcal{N}(0, R),\end{aligned}$$

the filter alternates a *prediction* step and an *update* step. The key equations are:

- **Predict:**

$$\hat{x}_{k|k-1} = A \hat{x}_{k-1|k-1}, \quad P_{k|k-1} = A P_{k-1|k-1} A^\top + Q.$$

- **Update:**

$$\begin{aligned}K_k &= P_{k|k-1} H^\top (H P_{k|k-1} H^\top + R)^{-1}, \quad \hat{x}_{k|k} = \hat{x}_{k|k-1} + K_k (y_k - H \hat{x}_{k|k-1}), \\P_{k|k} &= (I - K_k H) P_{k|k-1}.\end{aligned}$$

Here  $\hat{x}_{k|k-1}$  is the one-step-ahead state estimate and  $P$  its covariance. The **Kalman gain**  $K_k$  optimally weights the innovation  $y_k - H \hat{x}_{k|k-1}$  by the prior uncertainty. These equations assume *linear* dynamics and observation. In practice, we often cite a textbook or tutorial (e.g. Kalman 1960) for these formulas. (For example, Haykin notes that backpropagation can be derived as a “degenerate” Kalman update when simplifying assumptions are made <sup>1</sup>.)

## Extended Kalman Filter (EKF) – Nonlinear Models

When the system or observation is nonlinear, the **Extended Kalman Filter** linearizes the model around the current estimate <sup>2</sup>. Suppose the (nonlinear) state transition and observation are

$$x_k = f(x_{k-1}) + w_{k-1}, \quad y_k = h(x_k) + v_k,$$

with  $w \sim \mathcal{N}(0, Q)$ ,  $v \sim \mathcal{N}(0, R)$ . The EKF approximates  $f$  and  $h$  by their first-order Taylor expansions at the mean estimate. Let

$$F_{k-1} = \left. \frac{\partial f}{\partial x} \right|_{x=\hat{x}_{k-1|k-1}}, \quad H_k = \left. \frac{\partial h}{\partial x} \right|_{x=\hat{x}_{k|k-1}}$$

be the Jacobians of  $f$  and  $h$  (computed at the latest estimates) <sup>2</sup> . Then the EKF updates are:

- **Predict:**

$$\hat{x}_{k|k-1} = f(\hat{x}_{k-1|k-1}), \quad P_{k|k-1} = F_{k-1} P_{k-1|k-1} F_{k-1}^\top + Q.$$

- **Update:**

$$\begin{aligned} K_k &= P_{k|k-1} H_k^\top (H_k P_{k|k-1} H_k^\top + R)^{-1}, \\ \hat{x}_{k|k} &= \hat{x}_{k|k-1} + K_k (y_k - h(\hat{x}_{k|k-1})), \\ P_{k|k} &= (I - K_k H_k) P_{k|k-1}. \end{aligned}$$

These mirror the linear-Kalman equations but with  $f, h$  in place of  $A, H$  and with linearization. The Jacobian  $H_k$  ensures that the gain “magnifies” only the part of the residual that affects the state <sup>3</sup> . (Welch & Bishop note that an EKF is simply a Kalman filter “that linearizes about the current mean” <sup>2</sup> .)

## EKF Applied to MLP Training

In training a multilayer perceptron (MLP), we treat the network weights as the hidden state to be estimated. That is, let  $\theta$  denote the vector of all network weights (and biases) flattened into one vector. This is the EKF *state*. Given an input  $x_t$  and network output  $h(\theta, x_t)$ , we form the observation model

$$y_t = h(\theta, x_t) + v_t,$$

where  $y_t$  is the target output and  $v_t$  is assumed Gaussian measurement noise. Thus the **observation function** is  $h(\theta, x)$ , the MLP’s forward pass. As Singhal & Wu (1989) and others describe, the state estimation viewpoint means “the state of the system being estimated is the weights in the network. The observations are the outputs of the network as training vectors are presented” <sup>4</sup> .

We typically assume a trivial *state transition*: the weights drift slowly. For example, one can use a **random-walk model**  $\theta_t = I \theta_{t-1} + \eta_{t-1}$  with small process noise  $\eta \sim \mathcal{N}(0, Q)$ . In this case  $f(\theta) = \theta$  and  $F_t = I$ . In practice  $Q$  is set small to allow slight adaptation (e.g.  $Q = 10^{-4} I$  in one example <sup>5</sup> ).

The **Jacobian of the observation** w.r.t. the weights is computed by differentiating the MLP’s outputs with respect to each weight. Concretely, for a single-output MLP,  $h(\theta, x)$  is scalar and  $H_t = \nabla_\theta h(\theta_{t|t-1}, x_t)$  is a row vector of partial derivatives. For a multi-output MLP,  $H_t$  is a matrix (stacking each output’s gradient). In effect,  $H_t$  is the network’s backprop gradient of output error w.r.t. each weight.

Putting this into the EKF formulas:

• **Prediction:** Since  $f(\theta) = \theta$ , we have

$$\hat{\theta}_{t|t-1} = \hat{\theta}_{t-1|t-1}, \quad P_{t|t-1} = P_{t-1|t-1} + Q.$$

(Here  $P$  is the weight covariance matrix.)

- **Update:** With new training sample  $(x_t, y_t)$ , compute the output  $h(\hat{\theta}_{t|t-1}, x_t)$  and the Jacobian  $H_t$ . Then

$$S_t = H_t P_{t|t-1} H_t^\top + R, \quad K_t = P_{t|t-1} H_t^\top S_t^{-1}.$$

The innovation is  $y_t - h(\hat{\theta}_{t|t-1}, x_t)$ . Update:

$$\begin{aligned} \hat{\theta}_{t|t} &= \hat{\theta}_{t|t-1} + K_t (y_t - h(\hat{\theta}_{t|t-1}, x_t)), \\ P_{t|t} &= (I - K_t H_t) P_{t|t-1}. \end{aligned}$$

In words, the weight vector is corrected by the scaled error, where the **Kalman gain**  $K_t$  adapts each weight's step size. High gain (small  $R$ , large  $P$ ) means trusting the data more. In practice one often uses diagonal  $R = \sigma^2 I$  matching output noise variance, and sets  $Q \ll R$  for stability. For example, Chang *et al.* set  $Q = 10^{-4} I$  and  $R = y_{\text{std}}^2 I$  in their EKF MLP demo <sup>5</sup>.

**Summary of mathematical correspondences:** The EKF treats the MLP's weight vector as the latent state. The forward pass  $h(\theta, x)$  is the **measurement function**. Its Jacobian  $H = \frac{\partial h}{\partial \theta}$  is computed by backprop. Then the standard EKF gain and update formulas (with  $f(\theta) = \theta$ ) give the weight update rule above. All symbols align:  $x \leftrightarrow \theta$ ,  $y \leftrightarrow$  network output,  $f = I$ ,  $h(\theta, x) =$  network output.

## Comparison with Backpropagation

Standard *backpropagation* uses gradient descent on the loss  $L(\theta)$ . For a squared-error loss  $\frac{1}{2} \|y_t - h(\theta, x_t)\|^2$ , one computes the gradient  $\nabla_\theta L = -(y_t - h(\theta, x_t)) \nabla_\theta h(\theta, x_t)$  and updates

$$\theta \leftarrow \theta - \eta \nabla_\theta L.$$

This update is similar in spirit to the EKF correction  $\Delta\theta = K_t(y_t - h)$ , but there are key differences:

- **Information Usage:** Backprop's step uses only the (scaled) gradient  $\nabla_\theta h$  and a fixed learning rate  $\eta$ . In contrast, the EKF gain  $K_t = P_{t|t-1} H_t^\top (H_t P_{t|t-1} H_t^\top + R)^{-1}$  automatically adapts each weight's effective step size using the estimated covariance  $P$ . In effect, EKF uses a second-order-like weighting (since  $P$  captures uncertainty/Hessian information). This richer information often leads to faster convergence. Indeed, Singhal & Wu and others found that EKF training often needs far fewer epochs than backprop <sup>6</sup>.
- **Covariance vs. Learning Rate:** Backprop uses a constant (or decayed) scalar step-size. EKF effectively uses a *matrix* gain  $K_t$ . This gain can be interpreted as a learned, data-adaptive learning rate for each weight, weighting directions by their uncertainty. In other words, EKF is a quasi-second-order method because  $P$  approximates (the inverse of) the Hessian of the loss surface.

- **Convergence and Stability:** EKF can converge in fewer iterations or handle curvature better, but each update is much more expensive (inverting an  $m \times m$  matrix where  $m$  is output dimension or, more generally, an update of the full covariance of size  $n \times n$  where  $n$  is number of weights). Backprop's cost is just one gradient pass. A practical tradeoff is that EKF may reach higher accuracy quickly (fewer epochs) at the cost of heavy per-epoch computation <sup>6</sup>. Moreover, as noted by Ruck *et al.*, if one makes certain simplifications EKF *degenerates* to backprop: backprop “discards a great deal of information” (essentially treating  $P$  as a scalar) <sup>1</sup>. In practice, EKF often attains better accuracy per epoch than plain backprop <sup>1</sup>, but it requires careful tuning of  $Q, R$  to remain stable (too small  $R$  or too large gain can cause divergence).

## Mapping to the JAX/Flax Implementation

The provided JAX/Flax code realizes this EKF in an automated way. The main points of correspondence are:

- **Parameter Flattening:** The network's parameters are stored as a PyTree (nested dict) in Flax. Before applying EKF, they are flattened into a single 1D vector so that the filter state is a vector  $\theta$ . The code likely uses JAX's `ravel_pytree`, which “ravel[s] a pytree of arrays down to a 1D array” and returns also a function to unflatten it <sup>7</sup>. In practice, one does:

```
flat_params, unflatten = jax.flatten_util.ravel_pytree(params)
```

to get a flat weight vector and a callable to reconstruct the PyTree. This ensures consistent ordering of weights and allows matrix operations on the state. (See JAX docs <sup>7</sup> for details.)

- **Jacobian Computation** (`_jacobian_impl`): To compute  $H_t = \nabla_{\theta} h(\theta, x)$ , the code uses JAX's autodiff. A typical approach is `jax.jacrev` (reverse-mode Jacobian) or `jax.jacfwd`. In context, `_jacobian_impl` likely wraps `jax.jacrev` over the apply function. For each data point, this returns the gradient of the MLP output(s) w.r.t. each flattened weight. This corresponds exactly to the EKF's measurement Jacobian.
- **Kalman Gain and Update** (`_update_impl`, `_batch_update`): The functions named `_update_impl` and `_batch_update` implement the EKF formula. Given predicted mean  $\theta_{t|t-1}$  and covariance  $P_{t|t-1}$ , and given observation  $y_t$  and predicted  $h(\theta_{t|t-1}, x_t)$  and Jacobian  $H_t$ , the code computes the **innovation covariance**  $S = H_t P H_t^T + R$  and the **Kalman gain**  $K = P H_t^T S^{-1}$ . It then updates the state mean  $\theta \leftarrow \theta + K(y_t - h)$  and covariance  $P \leftarrow (I - K H_t) P$ . In a batch update, these operations are vectorized or iterated over multiple samples, aggregating information. The code likely uses `jax.numpy` matrix ops (`np.linalg.solve` or `np.linalg.inv`) to compute  $S^{-1}$ . These steps directly implement the EKF math from above.
- **Parameter Reconstruction:** After the filter updates the flat vector  $\theta$ , the code must map it back to the model parameters. This is done via the `unflatten` function returned by `ravel_pytree`. The new `theta` is un-raveled to the original nested dict of weights, which can then be fed into the model. In Flax, one often uses `model.apply(unflatten(flat_theta), x)` to get predictions.

- **bfloat16 and Numerical Stability:** The code uses `bfloat16` dtype for computations of the EKF. The bfloat16 format trades precision for speed and memory: it has the same range as float32 (8 exponent bits) but only 7 bits of mantissa <sup>8</sup>. This is advantageous on TPUs or accelerators, as it “reduces the storage requirements and increase[s] the calculation speed of ML algorithms” <sup>8</sup>. However, lower precision can cause round-off error. The implementation must therefore handle numerical stability carefully. For example, the covariance  $P$  update  $(I - KH)P$  can introduce negative eigenvalues if done naively; some implementations use the Joseph form  $P \leftarrow (I - KH)P(I - KH)^T + K R K^T$  to preserve symmetry. Also, adding a small “jitter” (diagonal noise) to  $P$  or ensuring  $R$  is not too small can keep  $S$  well-conditioned. In practice, using `bfloat16` means that one might accumulate  $P$  and  $K$  in higher precision or clamp values to avoid NaNs. The exact code may cast to `float32` internally for matrix inversions or use double precision for  $P$  while storing in bfloat16. The net effect is that the EKF runs faster on hardware at the cost of some precision, a tradeoff noted in mixed-precision training literature (e.g. Google TPU guides <sup>9</sup> <sup>8</sup>).

**Summary of Code-Math Mapping:** The flat weight vector  $\theta$  and covariance  $P$  correspond to the EKF state. The function `_jacobian_impl` computes  $H_t = \partial h / \partial \theta$ . The functions `_update_impl` and `_batch_update` implement the Kalman gain  $K$ , state update  $\theta \leftarrow \theta + K(y - h)$ , and  $P \leftarrow (I - KH)P$ . The utility `ravel_pytree` handles flattening/unflattening of parameter PyTrees <sup>7</sup>. Using `bfloat16` (as in `jax.numpy.bfloat16`) accelerates computation by halving storage, while requiring care to maintain numerical stability <sup>8</sup>.

In summary, the EKF recasts MLP training as a nonlinear state estimation problem. We derive the EKF equations from the standard Kalman filter by linearizing the network outputs w.r.t. weights. This yields adaptive, data-driven weight updates that incorporate both the error and the sensitivity of the network (the Jacobian). Compared to plain gradient descent, EKF uses extra covariance information to potentially converge faster, at the cost of higher computation per update <sup>6</sup> <sup>1</sup>. The provided JAX/Flax code mirrors this math: it flattens parameters, auto-differentiates to get gradients (Jacobian), computes the Kalman gain and innovation, and updates the weights and covariance accordingly. These correspondences ensure the implementation faithfully realizes the EKF algorithm of Singhal & Wu (1989) within a modern deep-learning framework.

**Sources:** The above derivations follow the classic EKF development (e.g. Welch & Bishop 2001 <sup>2</sup>) and the description of applying EKF to MLPs in Singhal & Wu (1989) <sup>4</sup> <sup>1</sup>. We also rely on a recent JAX example and documentation to interpret the code functions <sup>10</sup> <sup>7</sup> <sup>8</sup>.

---

<sup>1</sup> <sup>4</sup> <sup>6</sup> [cse.fau.edu](https://www.cse.fau.edu/~xqzhu/courses/cap5615/reading/FeatureSelectionUsingAmultilayerperceptron.pdf)  
<https://www.cse.fau.edu/~xqzhu/courses/cap5615/reading/FeatureSelectionUsingAmultilayerperceptron.pdf>

<sup>2</sup> <sup>3</sup> [2 The Extended Kalman Filter \(EKF\)](https://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL_COPIES/WELCH/kalman.2.html)  
[https://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL\\_COPIES/WELCH/kalman.2.html](https://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL_COPIES/WELCH/kalman.2.html)

<sup>5</sup> <sup>10</sup> [Online learning for an MLP using extended Kalman filtering](https://probml.github.io/dynamax/notebooks/nonlinear_gaussian_ssm/ekf_mlp.html)  
[https://probml.github.io/dynamax/notebooks/nonlinear\\_gaussian\\_ssm/ekf\\_mlp.html](https://probml.github.io/dynamax/notebooks/nonlinear_gaussian_ssm/ekf_mlp.html)

<sup>7</sup> [jax.flatten\\_util.ravel\\_pytree — JAX documentation](https://docs.jax.dev/en/latest/_autosummary/jax.flatten_util.ravel_pytree.html)  
[https://docs.jax.dev/en/latest/\\_autosummary/jax.flatten\\_util.ravel\\_pytree.html](https://docs.jax.dev/en/latest/_autosummary/jax.flatten_util.ravel_pytree.html)

8 bfloat16 floating-point format - Wikipedia

[https://en.wikipedia.org/wiki/Bfloat16\\_floating-point\\_format](https://en.wikipedia.org/wiki/Bfloat16_floating-point_format)

9 Improve your model's performance with bfloat16 | Cloud TPU

<https://cloud.google.com/tpu/docs/bfloat16>