

# Quadratic Kalman

Cedric Brendel

December 12, 2025

## 1 Notation

The shorthand  $:= y$  is to be read as “is defined as”, e. g.  $x := 2$  is to be read as “ $x$  is defined as 2”.

By  $\mathcal{N}(\mu, \Sigma)$  we denote the normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$  and we write  $X \sim \mathcal{N}(\mu, \Sigma)$  if the random variable  $X$  is distributed as such. In case that  $\Sigma$  is positive definite, we denote by  $\|x\|_\Sigma = \sqrt{x^\top \Sigma^{-1} x}$  the *Mahalanobis norm* of  $x$ . If  $X$  is a random variable and  $\mathcal{F}$  is some information (e. g. the information  $\{Y = y\}$  for some second random variable  $Y$  and a realization  $y$  thereof), we denote by  $X|\mathcal{F}$  the random variable obtained by conditioning on the information  $\mathcal{F}$ . Its density at  $x$  is denoted by  $p(X = x|\mathcal{F})$ . The expectation and covariance of the random variable  $X$  are denoted by  $\mathbb{E}[X]$  and  $\text{Cov}[X]$  and the notations  $\mathbb{E}[X|\mathcal{F}]$  and  $\text{Cov}[X|\mathcal{F}]$  do not(!) denote the conditional expectation and covariance but instead the expectation and covariance of the random variable  $X|\mathcal{F}$ .

The notations  $\nabla_x f$  and  $H_x f$  denote the gradient and Hessian of the function  $f$  at the point  $x$ . By  $f \propto g$  we denote the relation that the (real valued) functions  $f$  and  $g$  live on the same domain and are scalar multiples of each other, i. e. that there is some real  $c \neq 0$  with  $f = c \cdot g$ . Similarly, by  $f \frown g$  we denote the relation that the (real valued) functions  $f$  and  $g$  live on the same domain and are related by an (invertible) affine transformation, e. g.  $f \frown 2 \cdot f + 1$  but usually not  $f \frown \log f$ .

## 2 The Kalman Filter

Assume we have points in time  $t = 1, \dots, T$  and associated states  $X_t \sim \mathcal{N}(x_t, P_t)$ . Assume that the states evolve as  $X_t = F_t X_{t-1} + W_t$  where  $F_t$  is the *state transition matrix* and  $W_t \sim \mathcal{N}(0, Q_t)$  is the *state transition noise*. Assume further that the states  $X_t$  cannot be observed directly but instead only the variables  $Y_t = H_t X_t + V_t$  can be observed, where  $H_t$  is the *observation matrix* and  $V_t \sim \mathcal{N}(0, R_t)$  is the *observation noise*. Assume further that some initial state  $X_0 \sim \mathcal{N}(x_0, P_0)$  is given and that the initial state and all  $W_t$  and  $V_t$  are mutually independent. The *Kalman filter* estimates the (realized) states  $x_1, \dots, x_T$  from observations  $y_1, \dots, y_T$ . More precisely, (unbiased) estimates  $X_{t|t}$  of  $x_t$  and estimates  $P_{t|t}$  of the covariance of the error  $X_t - X_{t|t}$  are constructed for each  $t = 1, \dots, T$  such that each estimate only uses the information  $Y_k = y_k$  for  $k = 1, \dots, t$ .

The Kalman filter is well-viewed through a Bayesian lens. It consists of two *phases*. The *prediction phase* that estimates the posterior state  $X_t | \{Y_{\leq t-1} = y_{\leq t-1}\}$  at time  $t$  from observations up to and including  $t-1$ . The *update phase* incorporates the newly obtained observation  $y_t$  to estimate the posterior  $X_t | \{Y_{\leq t} = y_{\leq t}\}$ . As we will see, in both cases these posteriors are Gaussian and hence it is enough to estimate the means  $x_{t|t-1}$  and  $x_{t|t}$  as well as the covariances  $P_{t|t-1}$  and  $P_{t|t}$ .

Consider first the prediction phase. Using marginalization we find that

$$p(X_t = x_t | Y_{\leq t-1} = y_{\leq t-1}) = \int \underbrace{p(X_t = x_t | X_{t-1} = x_{t-1})}_{=\mathcal{N}(F_t x_{t-1}, Q_t)} \cdot \underbrace{p(X_{t-1} = x_{t-1} | Y_{\leq t-1} = y_{\leq t-1})}_{=\mathcal{N}(x_{t-1|t-1}, P_{t-1|t-1})} dx_{t-1}$$

since  $p(X_t = x_t | X_{t-1} = x_{t-1}, Y_{\leq t-1} = y_{\leq t-1}) = p(X_t = x_t | X_{t-1} = x_{t-1})$  as the next state  $X_t$  is completely determined by the current state  $x_{t-1}$  and the state transition noise  $W_t$  – all independent of  $\{Y_{\leq t-1} = y_{\leq t-1}\}$ . Thus,  $X_t | \{Y_{\leq t-1} = y_{\leq t-1}\}$  is again Gaussian. We find that

$$\begin{aligned} x_{t|t-1} &:= \mathbb{E}[F_t \cdot X_{t-1} + W_t | Y_{\leq t-1} = y_{\leq t-1}] \\ &= F_t \cdot \mathbb{E}[X_{t-1} | Y_{\leq t-1} = y_{\leq t-1}] + \mathbb{E}[W_t | Y_{\leq t-1} = y_{\leq t-1}] \\ &= F_t \cdot x_{t-1|t-1} + \mathbb{E}[W_t] \\ &= F_t \cdot x_{t-1|t-1} \end{aligned}$$

as  $W_t$  is independent of  $\{Y_{\leq t-1} = y_{\leq t-1}\}$ . For the covariance observe that

$$X_t - x_{t|t-1} = F_t \cdot (X_{t-1} - x_{t-1|t-1}) + W_t$$

and hence that

$$\begin{aligned} P_{t|t-1} &:= \text{Cov}[X_t | Y_{\leq t-1} = y_{\leq t-1}] \\ &= \mathbb{E}[(X_t - x_{t|t-1})(X_t - x_{t|t-1})^\top | Y_{\leq t-1} = y_{\leq t-1}] \\ &\stackrel{(*)}{=} F_t \cdot \mathbb{E}[(X_{t-1} - x_{t-1|t-1})(X_{t-1} - x_{t-1|t-1})^\top | Y_{\leq t-1} = y_{\leq t-1}] \cdot F_t^\top + \mathbb{E}[W_t W_t^\top] \\ &= F_t \text{Cov}[X_{t-1} | Y_{\leq t-1} = y_{\leq t-1}] F_t^\top + Q_t \\ &= F_t P_{t-1|t-1} F_t^\top + Q_t \end{aligned}$$

where  $(*)$  holds as the cross terms vanish and the final expectation is unconditional since  $W_t$  is independent of  $X_{t-1} - x_{t-1|t-1}$  and  $\{Y_{\leq t-1} = y_{\leq t-1}\}$  respectively. Thus we find

$$X_t | \{Y_{\leq t-1} = y_{\leq t-1}\} \sim \mathcal{N}(x_{t|t-1}, P_{t|t-1})$$

with  $x_{t|t-1} = F_t x_{t-1|t-1}$  and  $P_{t|t-1} = F_t P_{t-1|t-1} F_t^\top + Q_t$ .

We can now investigate  $X_t | \{Y_{\leq t} = y_{\leq t}\}$ . Applying Bayes' rule for densities, we find that

$$\begin{aligned} p(X_t = x_t | Y_{\leq t} = y_{\leq t}) &= p(X_t = x_t | Y_t = y_t, Y_{\leq t-1} = y_{\leq t-1}) \\ &\propto p(Y_t = y_t | X_t = x_t, Y_{\leq t-1} = y_{\leq t-1}) \cdot p(X_t = x_t | Y_{\leq t-1} = y_{\leq t-1}) \\ &= \underbrace{p(Y_t = y_t | X_t = x_t)}_{\mathcal{N}(H_t x_t, R_t)} \cdot \underbrace{p(X_t = x_t | Y_{\leq t-1} = y_{\leq t-1})}_{\mathcal{N}(x_{t|t-1}, P_{t|t-1})}, \end{aligned}$$

yielding that  $X_t | \{Y_{\leq t} = y_{\leq t}\}$  is again Gaussian. (Up to affine transformation) the negative log-likelihood is then

$$\begin{aligned} -\log p(X_t = x_t | Y_{\leq t} = y_{\leq t}) &= -\log p(Y_t = y_t | X_t = x_t) - \log p(X_t = x_t | Y_{\leq t-1} = y_{\leq t-1}) \\ &\sim \frac{1}{2}(y_t - H_t x_t)^\top R_t^{-1}(y_t - H_t x_t) + \frac{1}{2}(x_t - x_{t|t-1})^\top P_{t|t-1}^{-1}(x_t - x_{t|t-1}) \\ &= \frac{1}{2}\|y_t - H_t x_t\|_{R_t}^2 + \frac{1}{2}\|x_t - x_{t|t-1}\|_{P_{t|t-1}}^2. \end{aligned}$$

Recall that for a Gaussian  $X \sim \mathcal{N}(\mu, \Sigma)$  the negative log-likelihood is  $-\log p_X(x) \sim \frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)$  and hence the mean  $\mu$  is the unique minimizer of  $-\log p_X$  while the Hessian at  $\mu$  is the precision matrix  $\Sigma^{-1}$ . Thus the posterior mean respectively covariance of  $X_t | \{Y_{\leq t} = y_{\leq t}\}$  can be found as the minimizer respectively inverse of the Hessian (at the minimizer) of this quadratic function.

### 3 “Quadratic Kalman”

Assume we have points in time  $t = 1, \dots, T$  and are given a (time-varying, strictly convex) quadratic objective  $q_t$ , represented by its gradient  $\nabla_0 q_t$  and its (positive definite) Hessian  $H_0 q_t$  at 0, i.e.  $q_t(x) \sim \frac{1}{2}x^\top (H_0 q_t)x + (\nabla_0 q_t)^\top x$ . Further, assume that the state  $X_t$  is random, potentially satisfies equality

constraints  $A_t X_t = b_t$  where the  $A_t$  are matrices of full row rank and  $b_t$  are vectors, evolves as  $X_t = F_t X_{t-1} + W_t$  where  $F_t$  is the *state transition matrix* and  $W_t \sim \mathcal{N}(0, Q_t)$  is the *state transition noise*, and finally  $X_0 \sim \mathcal{N}(x_0, P_0)$  is given.

We seek to find an estimate  $x_{t|t}$  of the (realized, hidden) state  $x_t$  and its covariance  $P_{t|t}$  at time  $t$ , given information up to an including time  $t$ .

The prediction step of “Quadratic Kalman” is completely analogous to the one of the Kalman filter

The difference to the Kalman filter lies in the update step. Recall that the update step of the Kalman filter aims to uncover the distribution of the (Gaussian) random variable  $X_t | \{Y_{\leq t} = y_{\leq t}\}$  and does so by considering the conditional negative log likelihood

$$-\log p(X_t = x_t | Y_{\leq t} = y_{\leq t}) \sim \frac{1}{2} \|x_t - x_{t|t-1}\|_{P_{t|t-1}}^2 + \frac{1}{2} \|y_t - H_t x_t\|_{R_t}^2$$

which consists of the squared deviation  $\|x_t - x_{t|t-1}\|_{P_{t|t-1}}^2$  from the prior  $x_{t|t-1}$  and of the squared deviation  $\|y_t - H_t x_t\|_{R_t}^2$  of the observation  $y_t$  from the predicted observation  $H_t x_t$ . “Quadratic Kalman” instead pretends there is some “virtual” information  $\mathcal{F}_t$  uncovered up to time  $t$  and aims to uncover the distribution of the variable  $X_t | \mathcal{F}_t$ . It does so by assuming that the conditional log likelihood of this variable is given by

$$\tilde{q}_t(x_t) := -\log p(X_t = x_t | \mathcal{F}_t) := \frac{1}{2} \|x_t - x_{t|t-1}\|_{P_{t|t-1}}^2 + q_t(x_t),$$

where  $q_t$  is an arbitrary quadratic objective specified by the user. As this negative log likelihood is Gaussian, such a random variable  $X_t | \mathcal{F}_t$  must be Gaussian.

If one wants to incorporate the constraint  $A_t x_t = b_t$  into this probabilistic setting one has to adjust the estimated state covariance  $P_{t|t}$ . Indeed, the stricter the constraint (i. e. the smaller the dimension of the (affine) subspace of feasible  $x_t$  with  $A_t x_t = b_t$ ) the less uncertainty there is in our estimate of  $x_{t|t}$ . In the extreme case of a single feasible solution there is no(!) uncertainty! In “Quadratic Kalman” this is done by considering the family of (negative log) likelihoods

$$-\log p_\varepsilon(X_t = x_t | \mathcal{F}_t) := -\log p(X_t = x_t | \mathcal{F}_t) + \frac{1}{2\varepsilon} \|A_t x_t - b_t\|^2$$

and considering the limit of the corresponding estimates of mean and covariance (of an associated Gaussian random variable) as  $\varepsilon \rightarrow 0$ .

It is now not too hard to derive explicit analytic solutions of the desired state estimate  $x_{t|t}$  and state covariance  $P_{t|t}$  for both the unconstrained and constrained case using the lemmata 1, 2 and 3. For this, notice that for a (non-degenerate) Gaussian random variable, the mean and covariance are given by the unique minimizer and the inverse of the Hessian (at the minimizer) of the negative log likelihood. Indeed, if  $X \sim \mathcal{N}(\mu, \Sigma)$  then  $-\log p(X = x) = \frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)$  and the unique minimizer is  $\mu$  and the Hessian at the minimizer is the precision matrix  $\Sigma^{-1}$ .

For these lemmata we require the gradient and Hessian of the negative log likelihoods at 0. We calculate that

$$g_t := \nabla_0 \tilde{q}_t = \nabla_0 q_t - P_{t|t-1}^{-1} x_{t|t-1}$$

and

$$H_t := H_0 \tilde{q}_t = H_0 q_t + P_{t|t-1}^{-1}.$$

In the unconstrained case, lemma 1 yields that the updated state estimate is given by

$$x_{t|t} = -H_t^{-1} g_t$$

and it is clear that the updated state covariance estimate is

$$P_{t|t} = H_t^{-1}.$$

In the constrained case, lemma 2 yields that

$$x_{t|t} = -H_t^{-1}(b_t + A_t^\top \lambda^*)$$

with

$$\lambda^* = -(A_t H_t^{-1} A_t^\top)^{-1}(b_t + A_t H_t^{-1} g_t),$$

since  $A_t$  has full row rank. Finally, lemma 3 yields that the desired limit of inverse Hessians is given by

$$P_{t|t} = H_t^{-1} - H_t^{-1} A^\top (A_t H_t^{-1} A_t^\top)^{-1} A_t H_t^{-1}.$$

We obtain the “Quadratic Kalman” algorithm:

---

**Algorithm:** Quadratic-Kalman

---

**Require:** Initial state and state covariance estimate  $x_0$  and  $P_0$ .

**Require:** State transition and noise matrices  $F_t, Q_t$ .

**Require:** Hessians  $H_0 q_t$  and gradients  $\nabla_0 q_t$ .

**Require:** Maybe constraint data  $A_t, b_t$ .

```

 $x_{\text{upd}} \leftarrow x_0$ 
 $P_{\text{upd}} \leftarrow P_0$ 
for  $t = 1, \dots, T$  do

    > Prediction Step
     $x_{\text{pred}} \leftarrow F_t x_{\text{upd}}$ 
     $P_{\text{pred}} \leftarrow F_t P_{\text{upd}} F_t^\top + Q_t$ 

    > Update step
     $g \leftarrow \nabla_0 q_t - P_{\text{pred}}^{-1} x_{\text{pred}}$  > gradient of updated objective
     $H \leftarrow H_t q_t + P_{\text{pred}}^{-1}$  > Hessian of updated objective
    if constraint then
         $\lambda \leftarrow -(A_t H^{-1} A_t^\top)^{-1}(b_t + AH^{-1}g)$ 
         $x_{\text{upd}} \leftarrow -H^{-1}(b_t + A_t^\top \lambda)$ 
         $P_{\text{upd}} \leftarrow H^{-1} - H^{-1} A_t^\top (A_t H^{-1} A_t^\top)^{-1} A_t H^{-1}$ 
    else
         $x_{\text{upd}} \leftarrow -H^{-1}g$ 
         $P_{\text{upd}} \leftarrow H^{-1}$ 
    end if

end for

```

---

## A Appendix

**Lemma 1.** Let  $H \in \mathbf{R}^{n \times n}$  be symmetric positive definite and  $g \in \mathbf{R}^n$  be arbitrary. Then the quadratic objective

$$q(x) = \frac{1}{2}x^\top Hx + g^\top x$$

with  $x \in \mathbf{R}^n$  is strictly convex with unique minimizer given by

$$x^* = -H^{-1}g.$$

*Proof.* Observe that for any  $x \in \mathbf{R}^n$  we have  $\nabla_x q = Hx + g$  as  $H^\top = H$  by symmetry and hence  $H_x q = H$ . Thus, the hessian  $H_x q$  of  $q$  at any  $x \in \mathbf{R}^n$  is positive definite and hence  $q$  is strictly convex. Thus, there is unique minimizer  $x^*$  of  $q$  that satisfies  $\nabla_{x^*} q = 0$ . Clearly,  $\nabla_{x^*} q = Hx^* + g = 0$  is equivalent to  $x^* = -H^{-1}g$  by invertibility of  $H$ .  $\square$

**Lemma 2.** Consider again the setting of lemma 1. Let further  $A \in \mathbf{R}^{m \times n}$  have full row rank and  $b \in \mathbf{R}^m$  be arbitrary. The program

$$\begin{aligned} & \min \frac{1}{2} x^\top Hx + g^\top x \\ & \text{s.t. } Ax = b \end{aligned}$$

has a unique minimizer given by

$$x^* = -H^{-1}(g + A^\top \lambda^*)$$

where

$$\lambda^* = -(AH^{-1}A^\top)^{-1}(b + AH^{-1}g)$$

*Proof.* The existence of the unique minimizer  $x^*$  follows since  $q(x) = \frac{1}{2}x^\top Hx + g^\top x$  is strictly convex on the affine, non-empty (by full row rank of  $A$ ) subspace  $\{x \mid Ax = b\}$ . If  $x^*$  is a global (hence local) optimum, the KKT-conditions assert that there is a  $\lambda^* \in \mathbf{R}^m$  such that  $Ax^* = b$  and

$$\nabla_{x^*} q + A^\top \lambda^* = Hx^* + g + A^\top \lambda^* = 0.$$

Since  $H$  is invertible we find that

$$x^* = -H^{-1}(g + A^\top \lambda^*).$$

Substituting  $x^*$  in  $Ax^* = b$  we find that

$$b = -AH^{-1}(g + A^\top \lambda^*)$$

and hence

$$\lambda^* = -(AH^{-1}A^\top)^{-1}(b + AH^{-1}g)$$

as claimed, since  $AH^{-1}A^\top$  is invertible since  $H^{-1}$  is and  $A$  a full row rank.  $\square$

**Lemma 3.** Consider again the setting of lemma 1. Let  $\varepsilon > 0$  and consider the augmented quadratic

$$\tilde{q}_\varepsilon(x) = \frac{1}{2}x^\top Hx + g^\top x + \frac{1}{2\varepsilon}\|Ax - b\|_2^2$$

and denote by  $H_\varepsilon$  is Hessian at 0. Then

$$(H_0\tilde{q}_\varepsilon)^{-1} \rightarrow H^{-1} - H^{-1}A^\top(AH^{-1}A^\top)^{-1}AH^{-1}$$

as  $\varepsilon \rightarrow 0$ .

*Proof.* Observe that

$$\begin{aligned} \tilde{q}_\varepsilon(x) &= \frac{1}{2}x^\top Hx + g^\top x + \frac{1}{2\varepsilon}\|Ax - b\|_2^2 \\ &= \frac{1}{2}x^\top Hx + g^\top x + \frac{1}{2\varepsilon}(Ax - b)^\top(Ax - b) \\ &= \frac{1}{2}x^\top Hx + g^\top x + \frac{1}{2\varepsilon}x^\top A^\top Ax + \frac{1}{\varepsilon}b^\top Ax + \frac{1}{2\varepsilon}b^\top b \\ &= \frac{1}{2}x^\top \left(H + \frac{1}{\varepsilon}A^\top A\right)x + \left(g - \frac{1}{\varepsilon}A^\top b\right)^\top x + \frac{1}{2\varepsilon}b^\top b \end{aligned}$$

and hence that  $H_0\tilde{q}_\varepsilon = H + \frac{1}{\varepsilon}A^\top A$ . The Woodbury matrix identity thus gives that

$$\begin{aligned} (H_0\tilde{q}_\varepsilon)^{-1} &= \left(H + \frac{1}{\varepsilon}A^\top A\right)^{-1} \\ &= \left(H + A^\top \cdot \frac{1}{\varepsilon}E \cdot A\right)^{-1} \\ &= H^{-1} - H^{-1}A^\top(\varepsilon E + AH^{-1}A^\top)^{-1}AH^{-1} \end{aligned}$$

which clearly converges to

$$H^{-1} - H^{-1}A^\top(AH^{-1}A^\top)^{-1}AH^{-1}$$

for  $\varepsilon \rightarrow 0$ , as claimed.  $\square$