

# Lecture 2: From Linear Regression to Kalman Filter and Beyond

Simo Särkkä

January 15, 2020

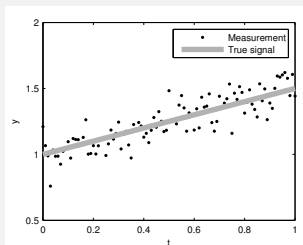
# Learning Outcomes

- 1 Summary of the Last Lecture
- 2 Batch and Recursive Estimation
- 3 Towards Bayesian Filtering
- 4 Kalman Filter and Bayesian Filtering and Smoothing
- 5 Summary

# Summary of the Last Lecture

- The purpose of is to estimate the **state of a time-varying system** from **noisy measurements** obtained from it.
- The **linear theory** dates back to **50's**, non-linear **Bayesian theory** was founded in **60's**.
- The efficient computational solutions can be divided into **prediction, filtering and smoothing**.
- **Applications**: tracking, navigation, telecommunications, audio processing, control systems, etc.
- The formal Bayesian estimation equations can be approximated by, e.g., **Gaussian approximations, Monte Carlo, or Gaussian mixtures**.
- **Formulating** physical systems as **state space models** is a challenging engineering topic as such.

# Batch Linear Regression [1/2]



- Consider the **linear regression model**

$$y_k = \theta_1 + \theta_2 t_k + \varepsilon_k, \quad k = 1, \dots, T,$$

with  $\varepsilon_k \sim \mathcal{N}(0, \sigma^2)$  and  $\theta = (\theta_1, \theta_2) \sim \mathcal{N}(\mathbf{m}_0, \mathbf{P}_0)$ .

- In **probabilistic notation** this is:

$$p(y_k | \theta) = \mathcal{N}(y_k | \mathbf{H}_k \theta, \sigma^2)$$

$$p(\theta) = \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{P}_0),$$

where  $\mathbf{H}_k = (1 \ t_k)$ .

- The **Bayesian batch solution** by the Bayes' rule:

$$\begin{aligned} p(\boldsymbol{\theta} | y_{1:T}) &\propto p(\boldsymbol{\theta}) \prod_{k=1}^T p(y_k | \boldsymbol{\theta}) \\ &= \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{P}_0) \prod_{k=1}^T \mathcal{N}(y_k | \mathbf{H}_k \boldsymbol{\theta}, \sigma^2). \end{aligned}$$

- The **posterior** is Gaussian

$$p(\boldsymbol{\theta} | y_{1:T}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_T, \mathbf{P}_T).$$

- The **mean and covariance** are given as

$$\begin{aligned} \mathbf{m}_T &= \left[ \mathbf{P}_0^{-1} + \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{H} \right]^{-1} \left[ \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{y} + \mathbf{P}_0^{-1} \mathbf{m}_0 \right] \\ \mathbf{P}_T &= \left[ \mathbf{P}_0^{-1} + \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{H} \right]^{-1}, \end{aligned}$$

where  $\mathbf{H}_k = (1 \ t_k)$ ,  $\mathbf{H} = (\mathbf{H}_1; \mathbf{H}_2; \dots; \mathbf{H}_T)$ ,  $\mathbf{y} = (y_1; \dots; y_T)$ .

- Assume that we have already computed the posterior distribution, which is **conditioned on the measurements up to  $k - 1$** :

$$p(\theta | y_{1:k-1}) = N(\theta | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}).$$

- Assume that we get the  **$k$ th measurement  $y_k$** . Using the equations from the previous slide we get

$$\begin{aligned} p(\theta | y_{1:k}) &\propto p(y_k | \theta) p(\theta | y_{1:k-1}) \\ &\propto N(\theta | \mathbf{m}_k, \mathbf{P}_k). \end{aligned}$$

- The **mean and covariance** are given as

$$\begin{aligned} \mathbf{m}_k &= \left[ \mathbf{P}_{k-1}^{-1} + \frac{1}{\sigma^2} \mathbf{H}_k^T \mathbf{H}_k \right]^{-1} \left[ \frac{1}{\sigma^2} \mathbf{H}_k^T y_k + \mathbf{P}_{k-1}^{-1} \mathbf{m}_{k-1} \right] \\ \mathbf{P}_k &= \left[ \mathbf{P}_{k-1}^{-1} + \frac{1}{\sigma^2} \mathbf{H}_k^T \mathbf{H}_k \right]^{-1}. \end{aligned}$$

- By the **matrix inversion lemma** (or Woodbury identity):

$$\mathbf{P}_k = \mathbf{P}_{k-1} - \mathbf{P}_{k-1} \mathbf{H}_k^T \left[ \mathbf{H}_k \mathbf{P}_{k-1} \mathbf{H}_k^T + \sigma^2 \right]^{-1} \mathbf{H}_k \mathbf{P}_{k-1}.$$

- Now the equations for the **mean and covariance** reduce to

$$\mathbf{S}_k = \mathbf{H}_k \mathbf{P}_{k-1} \mathbf{H}_k^T + \sigma^2$$

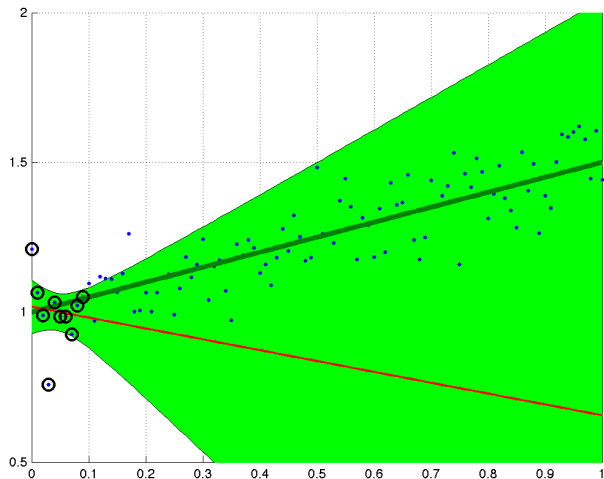
$$\mathbf{K}_k = \mathbf{P}_{k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1}$$

$$\mathbf{m}_k = \mathbf{m}_{k-1} + \mathbf{K}_k [y_k - \mathbf{H}_k \mathbf{m}_{k-1}]$$

$$\mathbf{P}_k = \mathbf{P}_{k-1} - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.$$

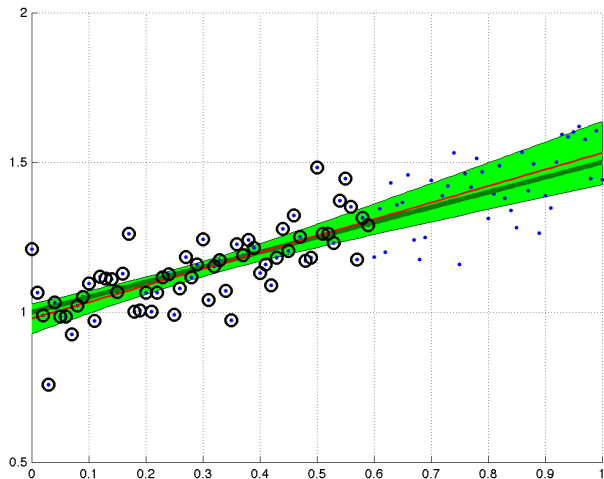
- Computing these for  $k = 0, \dots, T$  gives **exactly the linear regression solution**.
- A special case of **Kalman filter**.

# Recursive Linear Regression [3/4]

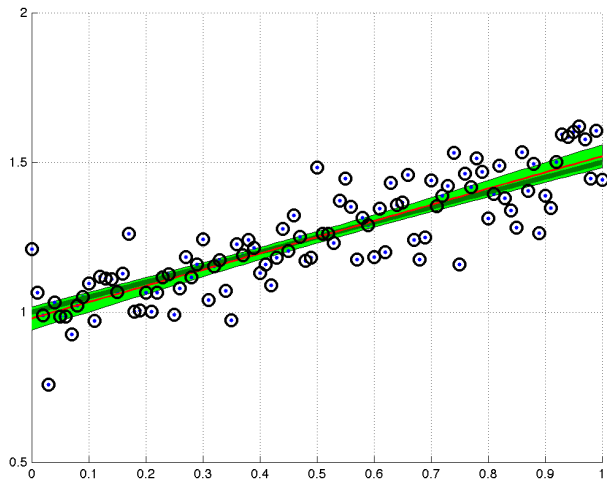




# Recursive Linear Regression [3/4]

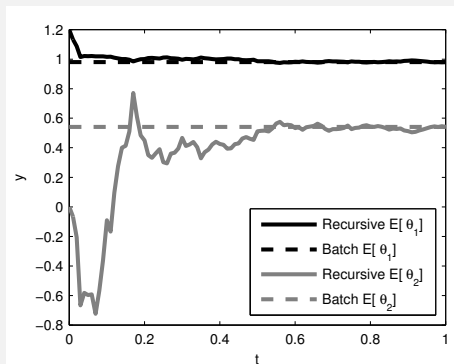


# Recursive Linear Regression [3/4]



# Recursive Linear Regression [4/4]

Convergence of the recursive solution to the batch solution – on the last step the solutions are exactly equal:



*General batch solution:*

- Specify the **measurement model**:

$$p(\mathbf{y}_{1:T} | \boldsymbol{\theta}) = \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}).$$

- Specify the **prior distribution**  $p(\boldsymbol{\theta})$ .
- Compute **posterior distribution** by the Bayes' rule:

$$p(\boldsymbol{\theta} | \mathbf{y}_{1:T}) = \frac{1}{Z} p(\boldsymbol{\theta}) \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}).$$

- Compute point estimates, moments, predictive quantities etc. from the posterior distribution.

# Batch vs. Recursive Estimation [2/2]

*General recursive solution:*

- Specify the **measurement likelihood**  $p(\mathbf{y}_k | \theta)$ .
- Specify the **prior distribution**  $p(\theta)$ .
- Process measurements  $\mathbf{y}_1, \dots, \mathbf{y}_T$  **one at a time**, starting from the prior:

$$p(\theta | \mathbf{y}_1) = \frac{1}{Z_1} p(\mathbf{y}_1 | \theta) p(\theta)$$

$$p(\theta | \mathbf{y}_{1:2}) = \frac{1}{Z_2} p(\mathbf{y}_2 | \theta) p(\theta | \mathbf{y}_1)$$

$$p(\theta | \mathbf{y}_{1:3}) = \frac{1}{Z_3} p(\mathbf{y}_3 | \theta) p(\theta | \mathbf{y}_{1:2})$$

$\vdots$

$$p(\theta | \mathbf{y}_{1:T}) = \frac{1}{Z_T} p(\mathbf{y}_T | \theta) p(\theta | \mathbf{y}_{1:T-1}).$$

- The result at the last step is the **batch solution**.

# Advantages of Recursive Solution

- The recursive solution can be considered as the **online learning** solution to the Bayesian learning problem.
- **Batch** Bayesian inference is a **special case of recursive** Bayesian inference.
- The **parameter** can be modeled to **change** between the measurement steps  $\Rightarrow$  **basis of filtering theory**.

# Drift Model for Linear Regression [1/3]

- Let assume **Gaussian random walk** between the measurements in the linear regression model:

$$\begin{aligned}p(y_k | \theta_k) &= \text{N}(y_k | \mathbf{H}_k \theta_k, \sigma^2) \\p(\theta_k | \theta_{k-1}) &= \text{N}(\theta_k | \theta_{k-1}, \mathbf{Q}) \\p(\theta_0) &= \text{N}(\theta_0 | \mathbf{m}_0, \mathbf{P}_0).\end{aligned}$$

- Again, assume that we already know

$$p(\theta_{k-1} | y_{1:k-1}) = \text{N}(\theta_{k-1} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}).$$

- The **joint distribution** of  $\theta_k$  and  $\theta_{k-1}$  is (due to Markovianity of dynamics!):

$$p(\theta_k, \theta_{k-1} | y_{1:k-1}) = p(\theta_k | \theta_{k-1}) p(\theta_{k-1} | y_{1:k-1}).$$

- Integrating over  $\theta_{k-1}$  gives:

$$p(\theta_k | y_{1:k-1}) = \int p(\theta_k | \theta_{k-1}) p(\theta_{k-1} | y_{1:k-1}) d\theta_{k-1}.$$

- This equation for **Markov processes** is called the **Chapman-Kolmogorov equation**.
- Because the distributions are Gaussian, the **result is Gaussian**

$$p(\theta_k | y_{1:k-1}) = N(\theta_k | \mathbf{m}_k^-, \mathbf{P}_k^-),$$

where

$$\mathbf{m}_k^- = \mathbf{m}_{k-1}$$

$$\mathbf{P}_k^- = \mathbf{P}_{k-1} + \mathbf{Q}.$$



- As in the pure recursive estimation, we get

$$\begin{aligned} p(\boldsymbol{\theta}_k | y_{1:k}) &\propto p(y_k | \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k | y_{1:k-1}) \\ &\propto \mathbf{N}(\boldsymbol{\theta}_k | \mathbf{m}_k, \mathbf{P}_k). \end{aligned}$$

- After applying the matrix inversion lemma, **mean and covariance** can be written as

$$\begin{aligned} \mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \sigma^2 \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k [y_k - \mathbf{H}_k \mathbf{m}_k^-] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T. \end{aligned}$$

- Again, we have derived a special case of the **Kalman filter**.
- The **batch version** of this solution would be **much more complicated**.

# State Space Notation

- In the previous slide we formulated the model as

$$p(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}) = N(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}, \mathbf{Q})$$
$$p(y_k | \boldsymbol{\theta}_k) = N(y_k | \mathbf{H}_k \boldsymbol{\theta}_k, \sigma^2)$$

- But in **Kalman filtering and control theory** the vector of parameters  $\boldsymbol{\theta}_k$  is usually called “state” and denoted as  $\mathbf{x}_k$ .
- More standard **state space notation**:

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}) = N(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{Q})$$
$$p(y_k | \mathbf{x}_k) = N(y_k | \mathbf{H}_k \mathbf{x}_k, \sigma^2)$$

- Or equivalently

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \mathbf{q}_{k-1}$$
$$y_k = \mathbf{H}_k \mathbf{x}_k + r_k,$$

where  $\mathbf{q}_{k-1} \sim N(\mathbf{0}, \mathbf{Q})$ ,  $r_k \sim N(0, \sigma^2)$ .

- The **canonical Kalman filtering model** is

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}) = N(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}_{k-1})$$

$$p(\mathbf{y}_k | \mathbf{x}_k) = N(\mathbf{y}_k | \mathbf{H}_k \mathbf{x}_k, \mathbf{R}_k).$$

- More often, this model can be seen **in the form**

$$\mathbf{x}_k = \mathbf{A}_{k-1} \mathbf{x}_{k-1} + \mathbf{q}_{k-1}$$

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{x}_k + \mathbf{r}_k.$$

- The Kalman filter actually calculates the following distributions:

$$p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-)$$

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) = N(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k).$$

- **Prediction step** of the Kalman filter:

$$\mathbf{m}_k^- = \mathbf{A}_{k-1} \mathbf{m}_{k-1}$$
$$\mathbf{P}_k^- = \mathbf{A}_{k-1} \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1}.$$

- **Update step** of the Kalman filter:

$$\mathbf{S}_k = \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k$$
$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T \mathbf{S}_k^{-1}$$
$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k [\mathbf{y}_k - \mathbf{H}_k \mathbf{m}_k^-]$$
$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.$$

- These equations will be derived from the general Bayesian filtering equations in the next lecture.

- Generic **non-linear state space models**

$$\mathbf{x}_k = \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1})$$

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}_k, \mathbf{r}_k).$$

- Generic **Markov models**

$$\mathbf{x}_k \sim p(\mathbf{x}_k | \mathbf{x}_{k-1})$$

$$\mathbf{y}_k \sim p(\mathbf{y}_k | \mathbf{x}_k).$$

- **Continuous-discrete** state space models involving **stochastic differential equations**:

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, t) + \mathbf{w}(t)$$

$$\mathbf{y}_k \sim p(\mathbf{y}_k | \mathbf{x}(t_k)).$$

- **Non-linear** state space model with **unknown parameters**:

$$\mathbf{x}_k = \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}, \boldsymbol{\theta})$$

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}_k, \mathbf{r}_k, \boldsymbol{\theta}).$$

- **General Markovian** state space model with **unknown parameters**:

$$\mathbf{x}_k \sim p(\mathbf{x}_k | \mathbf{x}_{k-1}, \boldsymbol{\theta})$$

$$\mathbf{y}_k \sim p(\mathbf{y}_k | \mathbf{x}_k, \boldsymbol{\theta}).$$

- Parameter estimation will be considered **later** – for now, we will attempt to **estimate the state**.
- Why **Bayesian filtering and smoothing** then?

# Bayesian Filtering, Prediction and Smoothing

- In principle, we could just use the (batch) **Bayes' rule**

$$\begin{aligned} & p(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{y}_1, \dots, \mathbf{y}_T) \\ &= \frac{p(\mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_T) p(\mathbf{x}_1, \dots, \mathbf{x}_T)}{p(\mathbf{y}_1, \dots, \mathbf{y}_T)}, \end{aligned}$$

- **Curse of computational complexity:** complexity grows more than linearly with number of measurements (typically we have  $O(T^3)$ ).
- Hence, we concentrate on the following:

- **Filtering distributions:**

$$p(\mathbf{x}_k \mid \mathbf{y}_1, \dots, \mathbf{y}_k), \quad k = 1, \dots, T.$$

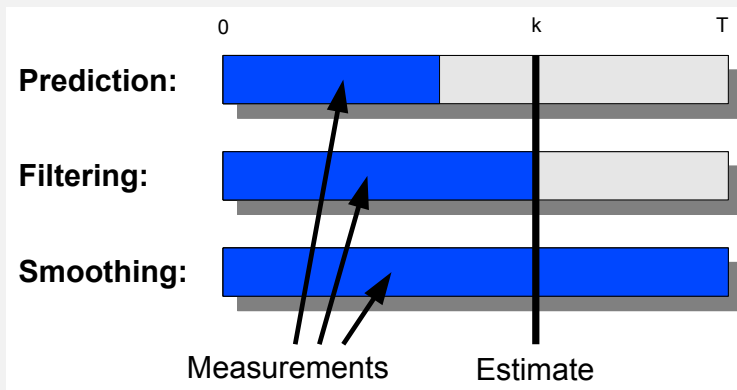
- **Prediction distributions:**

$$p(\mathbf{x}_{k+n} \mid \mathbf{y}_1, \dots, \mathbf{y}_k), \quad k = 1, \dots, T, \quad n = 1, 2, \dots,$$

- **Smoothing distributions:**

$$p(\mathbf{x}_k \mid \mathbf{y}_1, \dots, \mathbf{y}_T), \quad k = 1, \dots, T.$$

# Bayesian Filtering, Prediction and Smoothing (cont.)





- **Kalman filter** is the classical optimal filter for linear-Gaussian models.
- **Extended Kalman filter** (EKF) is linearization based extension of Kalman filter to non-linear models.
- **Unscented Kalman filter** (UKF) is sigma-point transformation based extension of Kalman filter.
- **Gauss-Hermite and Cubature Kalman filters** (GHKF/CKF) are numerical integration based extensions of Kalman filter.
- **Particle filter** forms a **Monte Carlo representation** (particle set) to the distribution of the state estimate.
- **Grid based filters** approximate the probability distributions on a finite grid.
- **Mixture Gaussian approximations** are used, for example, in multiple model Kalman filters and Rao-Blackwellized Particle filters.

- **Rauch-Tung-Striebel (RTS) smoother** is the closed form smoother for **linear Gaussian** models.
- **Extended, statistically linearized and unscented RTS smoothers** are the approximate nonlinear smoothers corresponding to EKF, SLF and UKF.
- **Gaussian RTS smoothers**: cubature RTS smoother, Gauss-Hermite RTS smoothers and various others
- **Particle smoothing** is based on approximating the smoothing solutions via **Monte Carlo**.
- **Rao-Blackwellized particle smoother** is a combination of particle smoothing and RTS smoothing.

# Summary

- **Linear regression problem** can be solved as **batch problem** or **recursively** – the latter solution is a special case of **Kalman filter**.
- A generic **Bayesian estimation problem** can also be solved as **batch problem** or **recursively**.
- If we let the linear regression **parameter change** between the measurements, we get a simple **linear state space model** – again solvable with **Kalman filtering model**.
- By **generalizing this idea** and the solution we get the **Kalman filter** algorithm.
- By further generalizing to **non-Gaussian models** results in generic **probabilistic state space models**.
- **Bayesian filtering and smoothing methods** solve Bayesian inference problems on state space models **recursively**.

[Linear regression with Kalman filter]