

PROBABILISTIC INFERENCE AND LEARNING

LECTURE 13

CONTINUOUS TIME SERIES

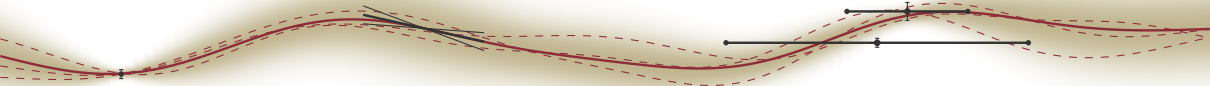
Philipp Hennig

28 November 2018

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

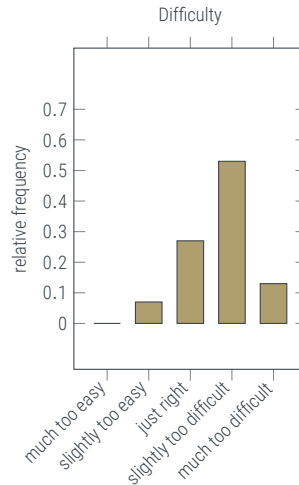
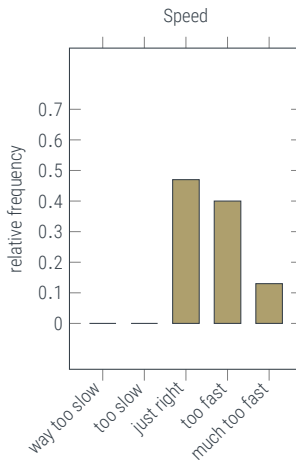
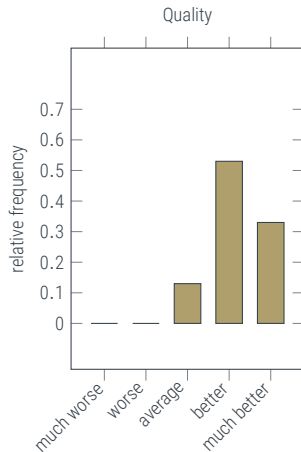


FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING



Last Lecture: Debrief

Feedback dashboard





Things you did not like:

- ✦ too fast
- ✦ exercises are again too hard

Things you did not understand:

- ✦ the Sobolev stuff
- ✦ do we have to understand the proofs?
- ✦ I don't know what Statistical Learning Theory is

Things you enjoyed:

- ✦ comparison of Bayesianism and Frequentism !!
- ✦ moral discussion at the end
- ✦ π -example

Overview of Lectures so far:

- | | |
|--|---|
| 0. Introduction to Reasoning under Uncertainty | 8. A practical GP example |
| 1. Probabilistic Reasoning | 9. Markov Chains, Time Series, Filtering |
| 2. Probabilities over Continuous Variables | 10. Classification |
| 3. Gaussian Probability Distributions | 11. Empirical Example of Classification |
| 4. Gaussian Parametric Regression | 12. Bayesianism and Frequentism |
| 5. More on Parametric Regression | 13. Stochastic Differential Equations (today) |
| 6. Gaussian Processes | 14. beyond Gaussians (WED) |
| 7. More on Kernels & GPs | |

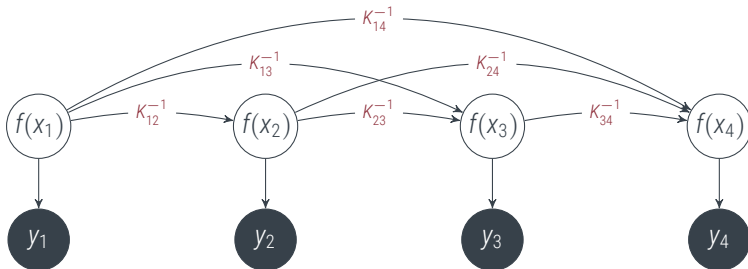
Today:

- ✦ Recap of Kalman Filters
- ✦ How to get from a Filter (A, Q) to a Gaussian process (m, k)

Graphical View I: “Full” GP

Recall from Lecture 1: Complexity of Inference can be controlled by Conditional Independence

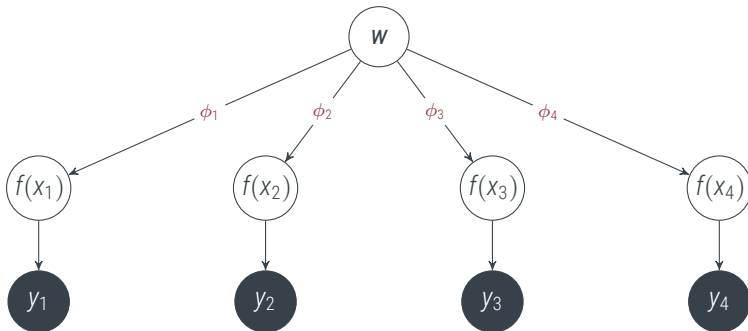
$$p(f) = \mathcal{GP}(f; 0, k) \quad p\left(\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{bmatrix}\right) = \mathcal{N}\left(0, \begin{bmatrix} K_{11}^{-1} & K_{12}^{-1} & K_{13}^{-1} & K_{14}^{-1} \\ & K_{22}^{-1} & K_{23}^{-1} & K_{24}^{-1} \\ & & K_{33}^{-1} & K_{34}^{-1} \\ & & & K_{44}^{-1} \end{bmatrix}^{-1}\right) \quad p(y | f) = \prod_i \mathcal{N}(y_i; f_i, \sigma^2)$$



Graphical View II: Parametric Model

Recall from Lecture 1: Complexity of Inference can be controlled by Conditional Independence

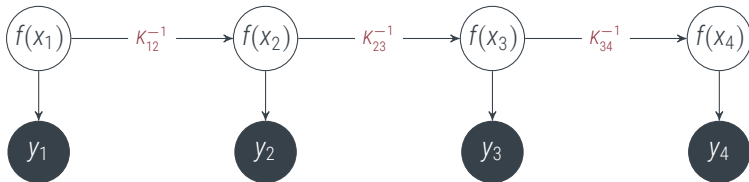
$$p(\mathbf{f}) = \mathcal{GP}(\mathbf{f}; 0, \Phi_X^T \Sigma \Phi_X) \quad p\left(\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{bmatrix} \middle| \mathbf{w}\right) = \prod_i \delta(f_i - \phi_i^T \mathbf{w}) \quad p(\mathbf{y} | \mathbf{f}) = \prod_i \mathcal{N}(y_i; f_i, \sigma^2)$$



Graphical View III: Markov Chain

Recall from Lecture 1: Complexity of Inference can be controlled by Conditional Independence

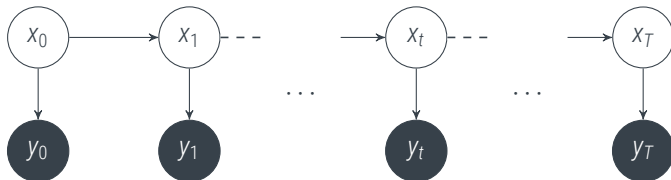
$$p(f) = \mathcal{GP}(f; 0, k) \quad p\left(\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{bmatrix}\right) = \mathcal{N}\left(0, \begin{bmatrix} K_{11}^{-1} & K_{12}^{-1} & 0 & 0 \\ K_{12}^{-1} & K_{22}^{-1} & K_{23}^{-1} & 0 \\ 0 & K_{23}^{-1} & K_{33}^{-1} & K_{34}^{-1} \\ 0 & 0 & K_{34}^{-1} & K_{44}^{-1} \end{bmatrix}^{-1}\right) \quad p(y | f) = \prod_i \mathcal{N}(y_i; f_i, \sigma^2)$$



Assume:

$$p(x_t | X_{0:t-1}) = p(x_t | x_{t-1})$$

and $p(y_t | X) = p(y_t | x_t)$



Filtering: $\mathcal{O}(T)$

predict:
$$p(x_t | Y_{0:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{0:t-1}) dx_{t-1} \quad (\text{Chapman-Kolmogorov Eq.})$$

update:
$$p(x_t | Y_{0:t}) = \frac{p(y_t | x_t) p(x_t | Y_{0:t-1})}{p(y_t)}$$

Smoothing: $\mathcal{O}(T)$

smooth:
$$p(x_t | Y) = p(x_t | Y_{0:t}) \int p(x_{t+1} | x_t) \frac{p(x_{t+1} | Y)}{p(x_{t+1} | Y_{1:t})} dx_{t+1}$$

Time Series:

- ✦ **Markov Chains** formalize the notion of a stochastic process with a *local finite memory*
- ✦ Inference over Markov Chains separates into three operations, that can be performed in *linear* time:

Filtering: $\mathcal{O}(T)$

predict:
$$p(x_t \mid Y_{0:t-1}) = \int p(x_t \mid x_{t-1}) p(x_{t-1} \mid Y_{0:t-1}) dx_{t-1} \quad (\text{Chapman-Kolmogorov Eq.})$$

update:
$$p(x_t \mid Y_{0:t}) = \frac{p(y_t \mid x_t) p(x_t \mid Y_{0:t-1})}{p(y_t)}$$

Smoothing: $\mathcal{O}(T)$

smooth:
$$p(x_t \mid Y) = p(x_t \mid Y_{0:t}) \int p(x_{t+1} \mid x_t) \frac{p(x_{t+1} \mid Y)}{p(x_{t+1} \mid Y_{1:t})} dx_{t+1}$$

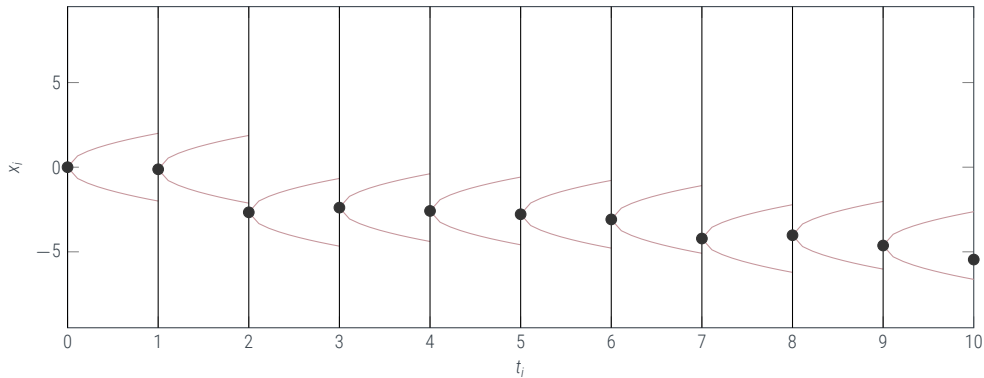
```

1 procedure INFERENCE( $Y, p(x_0), p(x_t | x_{t-1}) \forall t, p(y_t | x_t) \forall t$ )
2   for  $i=1, \dots, n$  do                                     // Filtering
3      $p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{0:t-1}) dx_{t-1}$  // Chapman-Kolmogorov eq.
4      $p(x_t | y_{1:t}) = p(y_t | x_t) p(x_t | Y_{0:t-1}) / p(y_t)$  // Update
5   end for
6   for  $i=n-1, \dots, 0$  do                                   // Smoothing
7      $p(x_t | Y) = p(x_t | Y_{0:t}) \int p(x_{t+1} | x_t) p(x_{t+1} | Y) p(x_{t+1} | Y_{1:t}) dx_{t+1}$ 
8   end for
9   return  $p(x_t | Y) \forall t = 0, \dots, n$  // return all marginals
10 end procedure

```



$$p(x(t_{i+1}) | x(t_i)) = \mathcal{N}(x_{i+1}; Ax_i, Q) \quad (\text{figure: } x_0 = 0, A = Q = 1)$$



$$p(x(t_i) \mid x(t_{i-1})) = \mathcal{N}(x_i; Ax_{i-1}, Q) \quad p(y_t \mid x_t) = \mathcal{N}(y_t; Hx_t, R)$$

predict:
$$p(x_t \mid Y_{1:t-1}) = \int p(x_t \mid x_{t-1}) p(x_{t-1} \mid Y_{0:t-1}) dx_{t-1} = \mathcal{N}(x_t, m_t^-, P_t^-)$$

$$= \mathcal{N}(x_t, Am_t, AP_tA^\top + Q)$$

update:
$$p(x_t \mid Y_{1:t}) = \frac{p(y_t \mid x_t) p(x_t \mid Y_{0:t-1})}{p(y_t)} = \mathcal{N}(x_t, m_t, P_t)$$

$$= \mathcal{N}(x_t, m_t^- + Kz, (I - KH)P_t^-)$$

$$K := P_t^- H^\top (HP_t^- H^\top + R)^{-1}, \quad z := y_t - Hm_t^-,$$

smooth:
$$p(x_t \mid Y) = p(x_t \mid Y_{0:t}) \int p(x_{t+1} \mid x_t) \frac{p(x_{t+1} \mid Y)}{p(x_{t+1} \mid Y_{1:t})} dx_{t+1} = \mathcal{N}(x_t, m_t^s, P_t^s)$$

$$= \mathcal{N}(x_t, m_t + G_t(m_{t+1}^s - m_{t+1}^-), P_t + G_t(P_{t+1}^s - P_{t+1}^-)G_t^\top)$$

$$G_t := P_t A^\top (P_{t+1}^-)^{-1}$$

$$p(x(t_i) \mid x(t_{i-1})) = \mathcal{N}(x_i; Ax_{i-1}, Q) \quad p(y_t \mid x_t) = \mathcal{N}(y_t; Hx_t, R)$$

(Kalman) Filter:

$m_t^- = Am_{t-1}$	predictive mean
$P_t^- = AP_{t-1}A^\top + Q$	predictive covariance
$z_t = y - Hm_t^-$	innovation residual
$S_t = HP_t^-H^\top + R$	innovation covariance
$K_t = P_t^-H^\top S_t^{-1}$	Kalman gain
$m_t = m_t^- + Kz_t$	estimation mean
$P_t = (I - KH)P_t^-$	estimation covariance

(Rauch Tung Striebel) Smoother:

$G_t = P_tA^\top(P_{t+1}^-)^{-1}$	RTS gain
$m_t^s = m_t + G_t(m_{t+1}^s - m_{t+1}^-)$	smoothed mean
$P_t^s = P_t + G_t(P_{t+1}^s - P_{t+1}^-)G_t^\top$	smoothed covariance



$$\delta t = 1 \quad Q_{\delta t} = 1$$



$$\delta t = 1/2 \quad Q_{\delta t} = 1/2$$



$$\delta t = 1/4 \quad Q_{\delta t} = \delta t$$

$$\delta t \rightarrow 0 \quad Q_{\delta t} = ???$$

For the limit $\delta t \rightarrow 0$ we would like to encode that $Q_{\delta t}/\delta t$ approaches some kind of finite object (like a derivative, but sample paths from this (the Wiener) process) are almost surely not differentiable. So we introduce a new object: $Q_{dt} := d\omega$, known as the *Wiener measure*. (Nb: This is a non-standard construction. $d\omega$ can be defined more elegantly; but this goes beyond the scope of this course.)

For our purposes the (linear, time-invariant) **Stochastic Differential Equation (SDE)**

$$dx(t) = Fx(t) dt + L d\omega_t,$$

together with $x(t_0) = x_0$, describes the local behaviour of the (unique) Gaussian process with

$$\mathbb{E}(x(t)) =: m(t) = e^{F(t-t_0)}x_0 \quad \text{cov}(x(t_a), x(t_b)) =: k(t_a, t_b) = \int_{t_0}^{\min t_a, t_b} e^{F(t_a-\tau)} LL^\top e^{F^\top(t_b-\tau)} d\tau$$

This GP is known as the **solution** of the SDE. It gives rise to the discrete-time stochastic recurrence relation $p(x_{t_{i+1}} | x_{t_i}) = \mathcal{N}(x_{t_{i+1}}; A_{t_i}x_{t_i}, Q_{t_i})$ with

$$A_{t_i} = e^{F(t_{i+1}-t_i)} \quad \text{and} \quad Q_{t_i} = \int_0^{t_{i+1}-t_i} e^{F\tau} LL^\top e^{F^\top\tau} d\tau.$$

Matrix exponential: $e^X := \sum_{i=0}^{\infty} \frac{X^i}{i!}$. Thus : $e^0 = I$, $(e^X)^{-1} = e^{-X}$, $X = VDV^{-1} \Rightarrow Ve^D V^{-1}$, $e^{\text{diag}_j d_j} = \text{diag}_j e^{d_j}$, $\det e^X = e^{\text{tr} X}$.

$$dx(t) = Fx(t) dt + L d\omega_t$$

$$\mathbb{E}(x(t)) =: m(t) = e^{F(t-t_0)}x_0 \quad \text{cov}(x(t_a), x(t_b)) =: k(t_a, t_b) = \int_{t_0}^{\min t_a, t_b} e^{F(t_a-\tau)} LL^\top e^{F^\top(t_b-\tau)} d\tau$$

$$A_{t_i} = e^{F(t_{i+1}-t_i)}$$

$$Q_{t_i} = \int_0^{t_{i+1}-t_i} e^{F\tau} LL^\top e^{F^\top\tau} d\tau$$



$$dx(t) = Fx(t) dt + L d\omega_t$$

$$\mathbb{E}(x(t)) =: m(t) = e^{F(t-t_0)}x_0 \quad \text{cov}(x(t_a), x(t_b)) =: k(t_a, t_b) = \int_{t_0}^{\min t_a, t_b} e^{F(t_a-\tau)} LL^\top e^{F^\top(t_b-\tau)} d\tau$$

$$A_{t_i} = e^{F(t_{i+1}-t_i)} \quad Q_{t_i} = \int_0^{t_{i+1}-t_i} e^{F\tau} LL^\top e^{F^\top\tau} d\tau$$

The scaled Wiener process

$$\begin{aligned} F = 0, L = \theta & \Rightarrow & m(t) &= x_0 & k(t_a, t_b) &= \theta^2(\min(t_a, t_b) - t_0) \\ & & A &= I & Q_{t_i} &= \theta^2(t_{i+1} - t_i) \end{aligned}$$



$$dx(t) = Fx(t) dt + L d\omega_t$$

$$\mathbb{E}(x(t)) =: m(t) = e^{F(t-t_0)}x_0 \quad \text{cov}(x(t_a), x(t_b)) =: k(t_a, t_b) = \int_{t_0}^{\min t_a, t_b} e^{F(t_a-\tau)} L L^\top e^{F^\top(t_b-\tau)} d\tau$$

$$A_{t_i} = e^{F(t_{i+1}-t_i)} \quad Q_{t_i} = \int_0^{t_{i+1}-t_i} e^{F\tau} L L^\top e^{F^\top\tau} d\tau$$

The Ornstein-Uhlenbeck process

$$F = -\frac{1}{\lambda}, L = \frac{2\theta}{\sqrt{\lambda}} \quad \Rightarrow \quad m(t) = x_0 e^{-\frac{t-t_0}{\lambda}} \quad k(t_a, t_b) = \theta^2 \left(e^{-\frac{|t_a-t_b|}{\lambda}} - e^{\frac{2t_0-t_a-t_b}{\lambda}} \right)$$

(cf. Exercises) $A = e^{-\delta_t/\lambda} \quad Q_{t_i} = \theta^2 \left(1 - e^{-2\delta_t/\lambda} \right)$

$$dx(t) = Fx(t) dt + L d\omega_t$$

- ✦ So far, we have seen examples with $x(t) \in \mathbb{R}$.
- ✦ But F and L can also be matrices. Consider the example

$$x = \begin{bmatrix} x_{(1)} \\ x_{(2)} \end{bmatrix} \quad F = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \quad L = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

That is:

$$\begin{bmatrix} dx_{(1)}(t) \\ dx_{(2)}(t) \end{bmatrix} = \begin{bmatrix} x_{(2)}(t) dt + 0 d\omega \\ 0 dt + d\omega \end{bmatrix} \Rightarrow x_{(1)}(t) = \int_{t_0}^t x_{(2)}(t) dt + [x_0]_1$$

- ✦ Finding A, Q is an exercise.

Inference doesn't get any faster than this

- Consider a *linear, time-invariant (LTI)* system (A, Q, H, R)

$$p(x_{i+1} \mid x_i) = \mathcal{N}(x_{i+1}; Ax_i, Q) \quad p(y_i \mid x_i) = \mathcal{N}(y_i; Hx_i, R)$$

```
1 procedure FILTER( $m_{t-1}, P_{t-1}, A, Q, H, R, y$ )
2    $m_t^- = Am_{t-1}$  // predictive mean
3    $P_t^- = AP_{t-1}A^\top + Q$  // predictive covariance
4    $z = y - Hm_t^-$  // residual
5    $S = HP_t^-H^\top + R$  // innovation covariance
6    $K = P_t^-H^\top S^{-1}$  // gain
7    $m_t = m_t^- + Kz$  // updated mean
8    $P_t = (I - KH)P_t^-$  // updated covariance
9   return  $(m_t, P_t), (m_t^-, P_t^-)$ 
10 end procedure
```

- P^- follows a **Discrete-time algebraic Riccati Equation (DARE)**

$$P_{i+1}^- = AP_i^-A^\top - AP_i^-H^\top(HP_i^-H^\top + R)^{-1}HP_i^-A^\top + Q = Q + A((P_i^-)^{-1} + H^\top R^{-1}H)^{-1}A^\top.$$

Inference doesn't get any faster than this

$$P_{i+1}^- = AP_i^- A^\top - AP_i^- H^\top (HP_i^- H^\top + R)^{-1} HP_i^- A^\top + Q = Q + A((P_i^-)^{-1} + H^\top R^{-1} H)^{-1} A^\top.$$

Theorem

Consider the (symplectic) matrix
$$Z = \begin{bmatrix} A^\top + H^\top R^{-1} H A^{-1} Q & -H^\top R^{-1} H A^{-1} \\ -A^{-1} Q & A^{-1} \end{bmatrix}.$$

If Z has no eigenvalues on the unit circle, then exactly half its eigenvalues are inside the unit circle, and the DARE above converges to a finite fix-point, known as the **steady-state** predictive covariance. To find it, consider the matrix $U \in \mathbb{C}^{2N \times 2N}$ of the eigenvectors of Z . Assume that they are sorted such that the first N columns of U correspond to the eigenvectors for eigenvalues inside the unit circle (and the other N columns to those outside the circle). Then separate U into square $N \times N$ sub-matrices as

$$U = \begin{bmatrix} U_1 & U_3 \\ U_2 & U_4 \end{bmatrix}.$$

The steady-state predictive covariance is given by $P_\infty^- = U_2 U_1^{-1}$.

In steady-state, $P_i^- = P_\infty^-$ and $K = P_\infty^- H^\top (H P_\infty^- H^\top + R)^{-1}$, Filter turns into

$$m_t \leftarrow A m_{t-1} + K(y - H A m_{t-1})$$

★ Wiener process: $F = 0, L = \theta, \delta t = 1, R = \sigma^2, H = I \Rightarrow A = I, Q = \theta^2$

$$Z = \begin{bmatrix} A^\top + H^\top R^{-1} H A^{-1} Q & -H^\top R^{-1} H A^{-1} \\ -A^{-1} Q & A^{-1} \end{bmatrix} = \begin{bmatrix} 1 + \theta^2/\sigma^2 & -1/\sigma^2 \\ -\theta^2 & 1 \end{bmatrix}$$

$$U = \begin{bmatrix} \frac{\sqrt{4\sigma^2/\theta^2 + 1} - 1}{2\sigma^2} & \frac{-\sqrt{4\sigma^2/\theta^2 + 1} - 1}{2\sigma^2} \\ 1 & 1 \end{bmatrix}, \quad P_\infty^- = \frac{2\sigma^2}{\sqrt{4\sigma^2/\theta^2 + 1} - 1}$$

$$m_t \leftarrow (1 - K)m_{t-1} + Ky \quad \text{with} \quad K = \frac{P_\infty^-}{P_\infty^- + \sigma^2}$$

The most basic kind of estimation

Steady-State filtering yields some old favorites

In steady-state, $P_i^- = P_\infty^-$ and $K = P_\infty^- H^\top (H P_\infty^- H^\top + R)^{-1}$, Filter turns into

$$m_t \leftarrow A m_{t-1} + K(y - H A m_{t-1})$$

★ Wiener process: $F = 0, L = \theta, \delta t = 1, R = \sigma^2, H = I \Rightarrow A = I, Q = \theta^2$

$$Z = \begin{bmatrix} A^\top + H^\top R^{-1} H A^{-1} Q & -H^\top R^{-1} H A^{-1} \\ -A^{-1} Q & A^{-1} \end{bmatrix} = \begin{bmatrix} 1 + \theta^2/\sigma^2 & -1/\sigma^2 \\ -\theta^2 & 1 \end{bmatrix}$$

$$U = \begin{bmatrix} \frac{\sqrt{4\sigma^2/\theta^2 + 1} - 1}{2\sigma^2} & \frac{-\sqrt{4\sigma^2/\theta^2 + 1} - 1}{2\sigma^2} \\ 1 & 1 \end{bmatrix}, \quad P_\infty^- = \frac{2\sigma^2}{\sqrt{4\sigma^2/\theta^2 + 1} - 1}$$

$$m_t \leftarrow (1 - K)m_{t-1} + Ky \quad \text{with} \quad K = \frac{P_\infty^-}{P_\infty^- + \sigma^2}$$

The steady-state of the Wiener-process filter is the **running average**.

Summary:

Markov Chains capture **finite memory** of a time series through conditional independence

Gauss-Markov models map this state to linear algebra

Kalman filter is the name for the corresponding algorithm

SDEs (Stochastic Differential Equations) are the continuous-time limit of discrete-time stochastic recurrence relations (in particular, linear SDEs are the continuous-time generalization discrete-time linear Gaussian systems)

Complexity of all necessary operations is **linear**, $\mathcal{O}(N)$ in the number of datapoints (as opposed to $\mathcal{O}(N^3)$ for general GPs).
(Although not shown, this includes hyperparameter inference!)

Steady-State update rules amount to various forms of **running averages** (the covariance is constant, the mean update rule is a running average)

HMMs are the generalisation of Gauss-Markov models to non-Gaussian generative relationships

For more on Gaussian and *approximately Gaussian filters* see, e.g.

Simo Särkkä. *Bayesian Filtering and Smoothing* Cambridge University Press, 2013

https://users.aalto.fi/~ssarkka/pub/cup_book_online_20131111.pdf

Summary of Gaussian Models

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

Gaussians provide the linear algebra of inference

if all joints are Gaussian and all observations are linear, all posteriors are Gaussian

- products of Gaussians are Gaussians

$$\begin{aligned}\mathcal{N}(x; a, A) \mathcal{N}(x; b, B) \\ = \mathcal{N}(x; c, C) \mathcal{N}(a; b, A + B)\end{aligned}$$

$$C := (A^{-1} + B^{-1})^{-1} \quad c := C(A^{-1}a + B^{-1}b)$$

- linear projections of Gaussians are Gaussians

$$\begin{aligned}p(z) &= \mathcal{N}(z; \mu, \Sigma) \\ \Rightarrow p(Az) &= \mathcal{N}(Az, A\mu, A\Sigma A^T)\end{aligned}$$

- marginals of Gaussians are Gaussians

$$\int \mathcal{N} \left[\begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$

- (linear) conditionals of Gaussians are Gaussians

$$\begin{aligned}p(x | y) &= \frac{p(x, y)}{p(y)} \\ &= \mathcal{N} \left(x; \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \right)\end{aligned}$$

Bayesian inference becomes linear algebra

If $p(x) = \mathcal{N}(x; \mu, \Sigma)$ and $p(y | x) = \mathcal{N}(y; A^T x + b, \Lambda)$, then

$$p(B^T x + c | y) = \mathcal{N}[B^T x + c; B^T \mu + c + B^T \Sigma A (A^T \Sigma A + \Lambda)^{-1} (y - A^T \mu - b), B^T \Sigma B - B^T \Sigma A (A^T \Sigma A + \Lambda)^{-1} A^T \Sigma B]$$

$$\text{prior } p(w) = \mathcal{N}(w; \mu, \Sigma) \Rightarrow p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood } p(y | w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 l) = \mathcal{N}(y; f_x, \sigma^2 l)$$

$$\begin{aligned} \text{posterior on } w \quad p(w | y, \phi_x) &= \mathcal{N}(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} (y - \phi_x^\top \mu), \\ &= \mathcal{N}\left(w; (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1}\right) \end{aligned}$$

$$\begin{aligned} \text{posterior on } f \quad p(f_x | y, \phi_x) &= \mathcal{N}(f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} (y - \phi_x^\top \mu), \\ &\quad \phi_x^\top \Sigma \phi_x - \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} \phi_x^\top \Sigma \phi_x) \\ &= \mathcal{N}\left(f_x; \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \phi_x^\top\right) \end{aligned}$$

$$\text{model evidence } p(y | \phi, \mu, \Sigma) = \mathcal{N}(y; \phi_x^\top \mu, \phi_x^\top \Sigma \phi_x + \sigma^2 l)$$

The model (ϕ, μ, Σ) can be learnt by **hierarchical Bayesian inference** (i.e. using the model evidence as a likelihood). Since the likelihood is not linear Gaussian for ϕ , approximate inference has to be used, e.g. MAP inference.

$$\text{prior} \quad p(f) = \mathcal{GP}(f; m, k)$$

$$\text{likelihood} \quad p(y \mid w, \phi_X) = \mathcal{N}(y; f_X, \sigma^2 I)$$

$$\text{posterior} \quad p(f \mid y, \phi_X) = \mathcal{GP}(f_X; m_X + k_{XX}(k_{XX} + \sigma^2 I)^{-1}(y - m_X), \quad k_X X - k_{XX}(k_{XX} + \sigma^2 I)^{-1}k_{XX})$$

Definition (Gaussian process)

Let $\mu : \mathbb{X} \rightarrow \mathbb{R}$ be any function, $k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ be a Mercer kernel. A **Gaussian process** $p(f) = \mathcal{GP}(f; \mu, k)$ is a probability distribution over the function $f : \mathbb{X} \rightarrow \mathbb{R}$, such that every finite restriction to function values $f_X := [f_{x_1}, \dots, f_{x_N}]$ is a Gaussian distribution $p(f_X) = \mathcal{N}(f_X; \mu_X, k_{XX})$.



- ✦ If k_1, k_2 are kernels, ϕ any function, then so are
 - ✦ $\alpha \cdot k_1(a, b)$ for $\alpha \in \mathbb{R}_+$
 - ✦ $k_1(\phi(c), \phi(d))$ for $c, d \in \mathbb{Y}$
 - ✦ $k_1(a, b) + k_2(a, b)$
 - ✦ $k_1(a, b) \cdot k_2(a, b)$
- ✦ for any kernel k there exists a unique *reproducing kernel Hilbert space* (RKHS) \mathcal{H}_k , and vice versa.
- ✦ GPs are quite powerful: The posterior mean can approximate (“learn”) any function in the RKHS at a rate given by that of the posterior standard deviation.
- ✦ GPs are quite limited: If $f \notin \mathcal{H}_k$, they may converge **very** (e.g. logarithmically) slowly to the truth
- ✦ Gaussian process regression is closely related to **kernel ridge regression** (nonparametric least-squares):
 - ✦ the posterior mean is the kernel ridge / regularized kernel least-squares estimate in the RKHS \mathcal{H}_k .

$$m(x) = k_{xx}(k_{xx} + \sigma^2 I)^{-1} \mathbf{y} = \arg \min_{f \in \mathcal{H}_k} \|\mathbf{y} - f_x\|^2 + \|f\|_{\mathcal{H}_k}^2$$

- ✦ the posterior variance (**expected square error**) is the **worst-case square error** in the RKHS.

$$v(x) = k_{xx} - k_{xx}(k_{xx} + \sigma^2 I)^{-1} k_{xx} = \arg \max_{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \leq 1} \|f(x) - m(x)\|^2$$

Kalman filters allow constant-time estimation in time series models

$$p(x(t_i) \mid x(t_{i-1})) = \mathcal{N}(x_i; Ax_{i-1}, Q) \quad p(y_t \mid x_t) = \mathcal{N}(y_t; Hx_t, R)$$

predict:
$$p(x_t \mid Y_{1:t-1}) = \int p(x_t \mid x_{t-1}) p(x_{t-1} \mid Y_{0:t-1}) dx_{t-1} = \mathcal{N}(x_t, m_t^-, P_t^-)$$

$$= \mathcal{N}(x_t, Am_t, AP_tA^\top + Q)$$

update:
$$p(x_t \mid Y_{1:t}) = \frac{p(y_t \mid x_t) p(x_t \mid Y_{0:t-1})}{p(y_t)} = \mathcal{N}(x_t, m_t, P_t)$$

$$= \mathcal{N}(x_t, m_t^- + Kz, (I - KH)P_t^-)$$

$$K := P_t^- H^\top (HP_t^- H^\top + R)^{-1}, \quad z := y_t - Hm_t^-,$$

smooth:
$$p(x_t \mid Y) = p(x_t \mid Y_{0:t}) \int p(x_{t+1} \mid x_t) \frac{p(x_{t+1} \mid Y)}{p(x_{t+1} \mid Y_{1:t})} dx_{t+1} = \mathcal{N}(x_t, m_t^s, P_t^s)$$

$$= \mathcal{N}(x_t, m_t + G_t(m_{t+1}^s - m_{t+1}^-), P_t + G_t(P_{t+1}^s - P_{t+1}^-)G_t^\top)$$

$$G_t := P_t A^\top (P_{t+1}^-)^{-1}$$

- ★ Assume $p(f) = \mathcal{GP}(f; m, k)$, but non-Gaussian likelihood $p(y | f)$ (e.g. $= \sigma(y \cdot f)$)
- ★ Find maximum posterior probability for **latent f** at **training points**

$$\hat{f} = \arg \max \log p(\mathbf{f}_X | y)$$

- ★ Assign approximate Gaussian posterior at training points

$$q(\mathbf{f}_X) = \mathcal{N}(\mathbf{f}_X; \hat{f}, -(\nabla \nabla^\top \log p(\mathbf{f}_X | y)|_{\mathbf{f}_X = \hat{f}})^{-1}) =: \mathcal{N}(\mathbf{f}_X; \hat{f}, \hat{\Sigma})$$

- ★ approximate posterior **predictions** at f_x for **latent function**

$$\begin{aligned} q(f_x | y) &= \int p(f_x | \mathbf{f}_X) q(\mathbf{f}_X) d\mathbf{f}_X = \int \mathcal{N}(f_x; m_x + k_{xx} K_{XX}^{-1} (\mathbf{f}_X - m_X), k_{xx} - k_{xx} K_{XX}^{-1} k_{xx}) q(\mathbf{f}_X) d\mathbf{f}_X \\ &= \mathcal{N}(f_x; m_x + k_{xx} K_{XX}^{-1} (\hat{f} - m_X), k_{xx} - k_{xx} K_{XX}^{-1} k_{xx} + k_{xx} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{xx}) \end{aligned}$$

- ★ compute predictions for **label probabilities**:

$$\mathbb{E}_{p(f|y)}[\pi_x] \approx \mathbb{E}_q[\pi_x] = \int \sigma(f_x) q(f_x | y) df_x$$

Gaussian distributions play a fundamental role in probabilistic reasoning,
 similar to how linear functions play a fundamental role in algebra.
 Gaussian models provide a flexible and powerful modeling language.
 With approximate inference, they even scale to non-Gaussian observations.

But Gaussians do not solve *all* problems. Some open questions for the coming lectures:

- ✦ We've seen fully connected models and Markov Chains. Is there something in between? What is the *complexity* of inference?
- ✦ We've seen that even discrete observations can be described by latent Gaussians. But what if the *latent* space is discrete? How can we efficiently keep track of fundamentally dissimilar competing hypotheses?
- ✦ The numerical approximations we have encountered so far are cheap, but full of flaws.
 - ✦ **maximum a-posteriori** has no uncertainty!
 - ✦ **Laplace approximation** can be arbitrarily wrong!
 - ✦ **numerical integration** does not scale to many dimensions!

Can we do better?