

The Oxford Lectures

A Short Course in State Estimation

Timothy D. Barfoot

University of Toronto

Copyright © 2024

Course Outline

Lecture 1: Linear-Gaussian Estimation

Lecture 2: Nonlinearities and Lie Groups

Lecture 3: Continuous-Time Estimation

Lecture 4: Beyond MAP – Variational Inference

Lecture 1: Linear-Gaussian Estimation

A Short Course in State Estimation

Timothy D. Barfoot

University of Toronto

Copyright © 2024

Lecture Outline

Lecture 1: Linear-Gaussian Estimation

- Problem Setup

- Bayesian Inference

- Maximum A Posteriori

- Sparsity

- Existence and Uniqueness

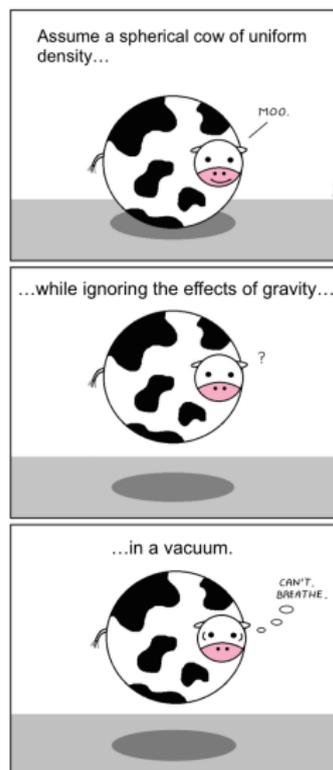
- Uncertainty

Lecture 2: Nonlinearities and Lie Groups

Lecture 3: Continuous-Time Estimation

Lecture 4: Beyond MAP – Variational Inference

Why linear-Gaussian?



- most of the classical results are for linear-Gaussian systems
- we can often find results without approximation
- we can get some intuition for real systems by studying them

“Equation (1.2-9) is a second order, nonlinear, vector, differential equation which has defied solution in its present form. It is here therefore we depart from the realities of nature to make some simplifying assumptions...”

– Bate et al., *Fundamentals of Astrodynamics* (1971)

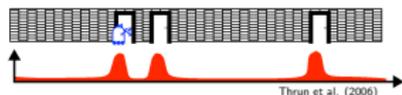
“Linear systems are very important because we can solve them exactly!”

– Sami Mikhail (heard third hand via Raja Mukherji)

What are the odds? (probability in one slide)

probability density function (PDF)

$$p(\mathbf{x})$$



conditional PDF (sensor model)

$$p(\mathbf{y}|\mathbf{x})$$

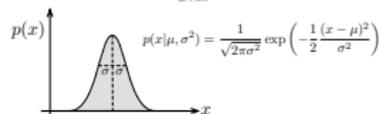
Bayes' rule (inference)

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$



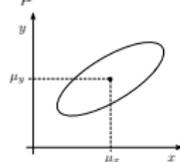
Gaussian PDF

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$



joint Gaussian PDF

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix}\right)$$

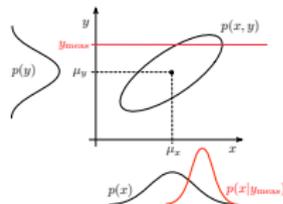


Gaussian marginalization

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$$

Gaussian conditioning

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx})$$



System

- we define our system using the following **linear, time-varying** models:

$$\text{motion model: } \mathbf{x}_k = \mathbf{A}_{k-1}\mathbf{x}_{k-1} + \mathbf{v}_k + \mathbf{w}_k, \quad k = 1 \dots K \quad (1a)$$

$$\text{observation model: } \mathbf{y}_k = \mathbf{C}_k\mathbf{x}_k + \mathbf{n}_k, \quad k = 0 \dots K \quad (1b)$$

where k is the discrete-time index and K its maximum

- the variables have the following meanings:

$$\text{system state : } \mathbf{x}_k \in \mathbb{R}^N$$

$$\text{initial state : } \mathbf{x}_0 \in \mathbb{R}^N \sim \mathcal{N}(\check{\mathbf{x}}_0, \check{\mathbf{P}}_0)$$

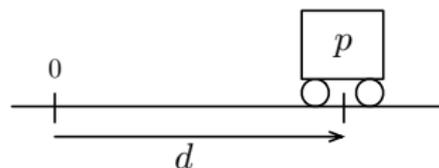
$$\text{input : } \mathbf{v}_k \in \mathbb{R}^N$$

$$\text{process noise : } \mathbf{w}_k \in \mathbb{R}^N \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k)$$

$$\text{measurement : } \mathbf{y}_k \in \mathbb{R}^M$$

$$\text{measurement noise : } \mathbf{n}_k \in \mathbb{R}^M \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$$

System example



$$d_k = p_k + n_k$$

$$\underbrace{d_k}_{\mathbf{y}_k} = \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_{\mathbf{C}_k} \underbrace{\begin{bmatrix} p_k \\ v_k \end{bmatrix}}_{\mathbf{x}_k} + \underbrace{n_k}_{\mathbf{n}_k}$$

observation model

$$\ddot{p}(t) = a(t) + w(t)$$

$$\dot{p}(t) = v(t)$$

$$\dot{v}(t) = a(t) + w(t)$$

$$p_k \approx p_{k-1} + T v_{k-1} + \frac{1}{2} T^2 (a_k + w_k)$$

$$v_k \approx v_{k-1} + T (a_k + w_k)$$

$$\underbrace{\begin{bmatrix} p_k \\ v_k \end{bmatrix}}_{\mathbf{x}_k} = \underbrace{\begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix}}_{\mathbf{A}_{k-1}} \underbrace{\begin{bmatrix} p_{k-1} \\ v_{k-1} \end{bmatrix}}_{\mathbf{x}_{k-1}} + \underbrace{\begin{bmatrix} \frac{1}{2} T^2 \\ T \end{bmatrix}}_{\mathbf{v}_k} a_k + \underbrace{\begin{bmatrix} \frac{1}{2} T^2 \\ T \end{bmatrix}}_{\mathbf{w}_k} w_k$$

motion model

What do we know?

- although we want to know the state of the system (at all times), we only have access to the following quantities, and must base our **estimate**, $\hat{\mathbf{x}}_k$, on just this information:
 - (i) the initial state knowledge, $\check{\mathbf{x}}_0$, and the associated covariance matrix, $\check{\mathbf{P}}_0$; sometimes we do not have this piece of information and must do without
 - (ii) the inputs, \mathbf{v}_k , which typically come from the output of our controller and so are known; we also have the associated process noise covariance, \mathbf{Q}_k
 - (iii) the measurements, $\mathbf{y}_{k,\text{meas}}$, which are **realizations** of the associated random variables, \mathbf{y}_k , and the associated covariance matrix, \mathbf{R}_k
- we will use $\hat{(\cdot)}$ to indicate **posterior** estimates (incorporating measurements) and $\check{(\cdot)}$ to indicate **prior** estimates (not incorporating measurements)

Problem statement

- we define the **state estimation problem** as follows:

Definition

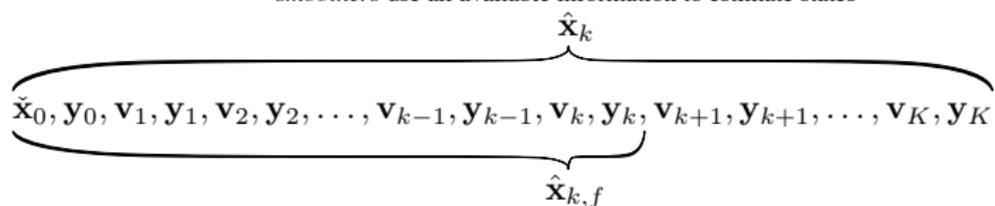
The problem of **state estimation** is to come up with an estimate, $\hat{\mathbf{x}}_k$, of the true *state* of a system, at one or more timesteps, k , given knowledge of the initial state, $\check{\mathbf{x}}_0$, a sequence of measurements, $\mathbf{y}_{0:K,\text{meas}}$, a sequence of inputs, $\mathbf{v}_{1:K}$, as well as knowledge of the system's motion and observation models.

- our approach will always be to attempt to come up with not only a state estimate, but also to quantify the **uncertainty** in that estimate

Batch is best

- we will begin by investigating **batch** linear-Gaussian techniques (sometimes called **smoothers**)
- the batch solution is very useful for computing state estimates after the fact because it uses all the measurements in the estimation of all the states at once (hence the usage of 'batch')
- a batch method cannot be used online since we cannot employ future measurements to estimate past states
- **recursive** state estimators (called **filters**) can be used online

smoothers use all available information to estimate states



filters only use past/current information to estimate states

All roads lead to Rome

- we will set up the batch **linear-Gaussian** estimation problem using two different paradigms:
 - (i) **Bayesian inference**; here we update a prior density over states (based on the initial state knowledge, inputs, and motion model) with our measurements, to produce a posterior (Gaussian) density over states
 - (ii) **maximum a posteriori**; here we employ optimization to find the most likely posterior state given the information we have (initial state knowledge, measurements, inputs)
- while these approaches are somewhat different in nature, it turns out that we arrive at the exact same answer for the **linear-Gaussian** problem
- this is because the full **Bayesian posterior is exactly Gaussian** and so the optimization approach will find the maximum (i.e., mode) of a Gaussian, and this is the same as the mean

A marriage of convenience

- we will combine the initial state knowledge and inputs in the following way:

$$\mathbf{v} = (\check{\mathbf{x}}_0, \mathbf{v}_{1:K}) = (\check{\mathbf{x}}_0, \mathbf{v}_1, \dots, \mathbf{v}_K) \quad (2)$$

since they are both related to our prior knowledge of the state

- we will also define

$$\mathbf{y} = \mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K) \quad (3)$$

to denote all of our measurements and

$$\mathbf{x} = \mathbf{x}_{0:K} = (\mathbf{x}_0, \dots, \mathbf{x}_K) \quad (4)$$

for our entire state (all timesteps)

Approach (i): Bayesian inference

- we would like to compute the **full Bayesian posterior**:

$$p(\mathbf{x}|\mathbf{y}, \mathbf{v}) = \frac{p(\mathbf{y}|\mathbf{x}, \mathbf{v})p(\mathbf{x}|\mathbf{v})}{p(\mathbf{y}|\mathbf{v})} \quad (5)$$

- we will take the two-step approach of building a joint density over states, inputs, measurements,

$$p(\mathbf{x}, \mathbf{y}|\mathbf{v}) = \underbrace{p(\mathbf{y}|\mathbf{x}, \mathbf{v})}_{\text{observations}} \underbrace{p(\mathbf{x}|\mathbf{v})}_{\text{prior}} \quad (6)$$

and then factor it the other way, keeping only the bit we want:

$$p(\mathbf{x}, \mathbf{y}|\mathbf{v}) = \underbrace{p(\mathbf{x}|\mathbf{y}, \mathbf{v})}_{\text{posterior}} \underbrace{p(\mathbf{y}|\mathbf{v})}_{\text{discard}} \quad (7)$$

Prior

- the **lifted mean** is then

$$\check{\mathbf{x}} = E[\mathbf{x}] = E[\mathbf{A}(\mathbf{v} + \mathbf{w})] = \mathbf{A}\mathbf{v} \quad (11)$$

- the **lifted covariance** is

$$\check{\mathbf{P}} = E [(\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^T] = \mathbf{A}\mathbf{Q}\mathbf{A}^T \quad (12)$$

where $\mathbf{Q} = E[\mathbf{w}\mathbf{w}^T] = \text{diag}(\check{\mathbf{P}}_0, \mathbf{Q}_1, \dots, \mathbf{Q}_K)$

- our **prior** can then be neatly expressed as

$$p(\mathbf{x}|\mathbf{v}) = \mathcal{N}(\check{\mathbf{x}}, \check{\mathbf{P}}) = \mathcal{N}(\mathbf{A}\mathbf{v}, \mathbf{A}\mathbf{Q}\mathbf{A}^T) \quad (13)$$

Measurements

- we now want to compute $p(\mathbf{x}, \mathbf{y} | \mathbf{v}) = p(\mathbf{y} | \mathbf{x}, \mathbf{v})p(\mathbf{x} | \mathbf{v})$
- we will use the measurement model

$$\mathbf{y}_k = \mathbf{C}_k \mathbf{x}_k + \mathbf{n}_k \quad (14)$$

- this can also be written in lifted form as

$$\mathbf{y} = \mathbf{C} \mathbf{x} + \mathbf{n} \quad (15)$$

where \mathbf{n} is the lifted form of the measurement noise and

$$\mathbf{C} = \text{diag}(\mathbf{C}_0, \mathbf{C}_1, \dots, \mathbf{C}_K) \quad (16)$$

is the lifted observation matrix

Joint Gaussian

- the joint likelihood of the prior lifted state and the measurements can now be written as

$$p(\mathbf{x}, \mathbf{y} | \mathbf{v}) = \mathcal{N} \left(\begin{bmatrix} \check{\mathbf{x}} \\ \mathbf{C}\check{\mathbf{x}} \end{bmatrix}, \begin{bmatrix} \check{\mathbf{P}} & \check{\mathbf{P}}\mathbf{C}^T \\ \mathbf{C}\check{\mathbf{P}} & \mathbf{C}\check{\mathbf{P}}\mathbf{C}^T + \mathbf{R} \end{bmatrix} \right) \quad (17)$$

where $\mathbf{R} = E[\mathbf{nn}^T] = \text{diag}(\mathbf{R}_0, \mathbf{R}_1, \dots, \mathbf{R}_K)$

- we can factor this according to

$$p(\mathbf{x}, \mathbf{y} | \mathbf{v}) = p(\mathbf{x} | \mathbf{y}, \mathbf{v}) p(\mathbf{y} | \mathbf{v}) \quad (18)$$

- we only care about the first factor, which is the **full Bayesian posterior**

Posterior

- this can be written, using the **Gaussian conditioning formula**, as

$$p(\mathbf{x}|\mathbf{v}, \mathbf{y}) = \mathcal{N}\left(\check{\mathbf{x}} + \check{\mathbf{P}}\mathbf{C}^T(\mathbf{C}\check{\mathbf{P}}\mathbf{C}^T + \mathbf{R})^{-1}(\mathbf{y} - \mathbf{C}\check{\mathbf{x}}), \right. \\ \left. \check{\mathbf{P}} - \check{\mathbf{P}}\mathbf{C}^T(\mathbf{C}\check{\mathbf{P}}\mathbf{C}^T + \mathbf{R})^{-1}\mathbf{C}\check{\mathbf{P}}\right) \quad (19)$$

- using the **Sherman-Morrison-Woodbury** identity, this can be manipulated into the following form:

$$p(\mathbf{x}|\mathbf{v}, \mathbf{y}) = \mathcal{N}\left(\underbrace{\left(\check{\mathbf{P}}^{-1} + \mathbf{C}^T\mathbf{R}^{-1}\mathbf{C}\right)^{-1} \left(\check{\mathbf{P}}^{-1}\check{\mathbf{x}} + \mathbf{C}^T\mathbf{R}^{-1}\mathbf{y}\right)}_{\hat{\mathbf{x}}, \text{ mean}}, \right. \\ \left. \underbrace{\left(\check{\mathbf{P}}^{-1} + \mathbf{C}^T\mathbf{R}^{-1}\mathbf{C}\right)^{-1}}_{\hat{\mathbf{P}}, \text{ covariance}}\right) \quad (20)$$

- we could use this as is, but there is a better way

Posterior

- we can rearrange the mean expression to arrive at

$$\underbrace{(\check{\mathbf{P}}^{-1} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C})}_{\hat{\mathbf{P}}^{-1}} \hat{\mathbf{x}} = \check{\mathbf{P}}^{-1} \check{\mathbf{x}} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{y} \quad (21)$$

and we see the inverse covariance appearing on the left-hand side

- substituting in $\check{\mathbf{x}} = \mathbf{A} \mathbf{v}$ and $\check{\mathbf{P}}^{-1} = (\mathbf{A} \mathbf{Q} \mathbf{A}^T)^{-1} = \mathbf{A}^{-T} \mathbf{Q}^{-1} \mathbf{A}^{-1}$ we can rewrite this as

$$\underbrace{(\mathbf{A}^{-T} \mathbf{Q}^{-1} \mathbf{A}^{-1} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C})}_{\hat{\mathbf{P}}^{-1}} \hat{\mathbf{x}} = \mathbf{A}^{-T} \mathbf{Q}^{-1} \mathbf{v} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{y} \quad (22)$$

- or more compactly,

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \hat{\mathbf{x}} = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z} \quad (23)$$

$$\mathbf{z} = \begin{bmatrix} \mathbf{v} \\ \mathbf{y} \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \mathbf{A}^{-1} \\ \mathbf{C} \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{Q} & \\ & \mathbf{R} \end{bmatrix} \quad (24)$$

Approach (ii): maximum a posteriori (MAP)

- instead of computing the full Bayesian posterior, what if we just tried to find the **most likely state** of our system given the initial state knowledge, inputs, and measurements?
- to accomplish this, we could try to solve the following optimization problem:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}, \mathbf{v}) \quad (25)$$

which is to say that we want to find the best single estimate for the state of the system (at all timesteps), $\hat{\mathbf{x}}$, given the prior information, \mathbf{v} , and measurements, \mathbf{y}

MAP

- we begin by rewriting the MAP estimate using Bayes' rule:

$$\begin{aligned}\hat{\mathbf{x}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}, \mathbf{v}) = \arg \max_{\mathbf{x}} \frac{p(\mathbf{y}|\mathbf{x}, \mathbf{v})p(\mathbf{x}|\mathbf{v})}{p(\mathbf{y}|\mathbf{v})} \\ &= \arg \max_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\mathbf{v}) \quad (26)\end{aligned}$$

where we drop the denominator because it does not depend on \mathbf{x}

- we also drop \mathbf{v} in $p(\mathbf{y}|\mathbf{x}, \mathbf{v})$ since it does not affect \mathbf{y} in our system if \mathbf{x} is known (see observation model)

MAP

- if we assume that all of the noise variables, \mathbf{w}_k and \mathbf{n}_k for $k = 0 \dots K$, are uncorrelated, we can use Bayes' rule to factor $p(\mathbf{y}|\mathbf{x})$ in the following way:

$$p(\mathbf{y}|\mathbf{x}) = \prod_{k=0}^K p(\mathbf{y}_k | \mathbf{x}_k) \quad (27)$$

- furthermore, Bayes' rule allows us to factor $p(\mathbf{x}|\mathbf{v})$ as

$$p(\mathbf{x}|\mathbf{v}) = p(\mathbf{x}_0 | \check{\mathbf{x}}_0) \prod_{k=1}^K p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{v}_k) \quad (28)$$

Negative log likelihood

- to make the optimization easier, we insert the (negative) logarithm function, which is monotonically decreasing

$$\begin{aligned}\hat{\mathbf{x}} &= \arg \min_{\mathbf{x}} (-\ln(p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\mathbf{v}))) \\ &= \arg \min_{\mathbf{x}} \left(-\ln p(\mathbf{x}_0 | \check{\mathbf{x}}_0) - \sum_{k=1}^K \ln p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{v}_k) \right. \\ &\quad \left. - \sum_{k=0}^K \ln p(\mathbf{y}_k | \mathbf{x}_k) \right) \quad (29)\end{aligned}$$

- note, we are now minimizing due to the negative sign, which is more common in optimization

– the component Gaussian densities are given by

$$p(\mathbf{x}_0 | \check{\mathbf{x}}_0) = \frac{1}{\sqrt{(2\pi)^N \det \check{\mathbf{P}}_0}} \exp \left(-\frac{1}{2} (\mathbf{x}_0 - \check{\mathbf{x}}_0)^T \check{\mathbf{P}}_0^{-1} (\mathbf{x}_0 - \check{\mathbf{x}}_0) \right) \quad (30a)$$

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{v}_k) = \frac{1}{\sqrt{(2\pi)^N \det \mathbf{Q}_k}} \exp \left(-\frac{1}{2} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k)^T \right. \\ \left. \times \mathbf{Q}_k^{-1} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k) \right) \quad (30b)$$

$$p(\mathbf{y}_k | \mathbf{x}_k) = \frac{1}{\sqrt{(2\pi)^M \det \mathbf{R}_k}} \exp \left(-\frac{1}{2} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k)^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k) \right) \quad (30c)$$

– the component negative-log densities are given by

$$-\ln p(\mathbf{x}_0 | \check{\mathbf{x}}_0) = \frac{1}{2} (\mathbf{x}_0 - \check{\mathbf{x}}_0)^T \check{\mathbf{P}}_0^{-1} (\mathbf{x}_0 - \check{\mathbf{x}}_0) + \underbrace{\frac{1}{2} \ln \left((2\pi)^N \det \check{\mathbf{P}}_0 \right)}_{\text{independent of } \mathbf{x}} \quad (31a)$$

$$-\ln p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{v}_k) = \frac{1}{2} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k)^T \mathbf{Q}_k^{-1} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k) + \underbrace{\frac{1}{2} \ln \left((2\pi)^N \det \mathbf{Q}_k \right)}_{\text{independent of } \mathbf{x}} \quad (31b)$$

$$-\ln p(\mathbf{y}_k | \mathbf{x}_k) = \frac{1}{2} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k)^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k) + \underbrace{\frac{1}{2} \ln \left((2\pi)^M \det \mathbf{R}_k \right)}_{\text{independent of } \mathbf{x}} \quad (31c)$$

MAP as optimization

- from an optimization perspective, we want to solve the following problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) \quad (32)$$

- our **cost function** is

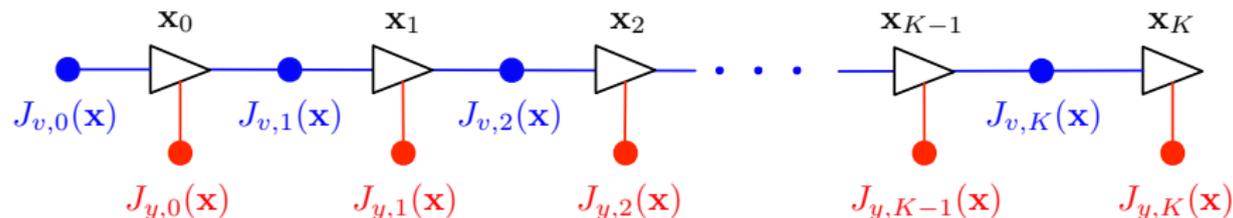
$$J(\mathbf{x}) = \sum_{k=0}^K (J_{v,k}(\mathbf{x}) + J_{y,k}(\mathbf{x})) \quad (33)$$

with

$$J_{v,k}(\mathbf{x}) = \begin{cases} \frac{1}{2} (\mathbf{x}_0 - \check{\mathbf{x}}_0)^T \check{\mathbf{P}}_0^{-1} (\mathbf{x}_0 - \check{\mathbf{x}}_0) & k = 0 \\ \frac{1}{2} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k)^T \\ \quad \times \mathbf{Q}_k^{-1} (\mathbf{x}_k - \mathbf{A}_{k-1} \mathbf{x}_{k-1} - \mathbf{v}_k) & k = 1 \dots K \end{cases} \quad (34a)$$

$$J_{y,k}(\mathbf{x}) = \frac{1}{2} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k)^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_k) \quad k = 0 \dots K \quad (34b)$$

MAP as a factor graph



$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) = \arg \min_{\mathbf{x}} \sum_{k=0}^K \left(J_{v,k}(\mathbf{x}) + J_{y,k}(\mathbf{x}) \right)$$

motion prior measurements

MAP as optimization

- we can express our cost function more compactly as

$$J(\mathbf{x}) = \frac{1}{2} (\mathbf{z} - \mathbf{H}\mathbf{x})^T \mathbf{W}^{-1} (\mathbf{z} - \mathbf{H}\mathbf{x}) \quad (35)$$

recalling our definitions from the full Bayesian approach:

$$\mathbf{z} = \begin{bmatrix} \mathbf{v} \\ \mathbf{y} \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \mathbf{A}^{-1} \\ \mathbf{C} \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{Q} & \\ & \mathbf{R} \end{bmatrix} \quad (36)$$

- we see that $J(\mathbf{x})$ is **exactly quadratic** in \mathbf{x}

MAP as optimization

- since $J(\mathbf{x})$ is exactly a paraboloid, we can find its minimum in closed form
- simply set the partial derivative with respect to the design variable, \mathbf{x} , to zero:

$$\left. \frac{\partial J(\mathbf{x})}{\partial \mathbf{x}^T} \right|_{\hat{\mathbf{x}}} = -\mathbf{H}^T \mathbf{W}^{-1} (\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}) = \mathbf{0} \quad (37a)$$

- rearranging we have

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \hat{\mathbf{x}} = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z} \quad (37b)$$

which is the identical linear system we arrived at in the full Bayesian approach

What gives?

- why do the full Bayesian and MAP approaches result in the same linear system if they are trying to do different things?

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \hat{\mathbf{x}} = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z} \quad (38)$$

- the explanation lies in the fact that the full Bayesian posterior is **exactly Gaussian** and the mean and mode (i.e., maximum) of a Gaussian are one and the same
- for nonlinear systems, these two approaches will diverge, because the posterior will not be exactly Gaussian

Cholesky solution

- once we have the **Cholesky factorization**, the solution of our linear system is straightforward

- **forward pass**: solve

$$\mathbf{L}\mathbf{d} = \mathbf{H}^T\mathbf{W}^{-1}\mathbf{z} \quad (44)$$

for \mathbf{d}

- **backward pass**: solve

$$\mathbf{L}^T\hat{\mathbf{x}} = \mathbf{d} \quad (45)$$

for $\hat{\mathbf{x}}$

- both passes can be done in $O(N^3(K + 1))$ time through forward/backward substitution owing to the sparse lower-triangular form of \mathbf{L}
- thus, the batch equations can be solved in computation time that scales linearly with the size of the state

Cholesky solution

- at the block level, the **forward pass**, $k = 1 \dots K$, is

$$\mathbf{L}_{k-1} \mathbf{L}_{k-1}^T = \mathbf{I}_{k-1} + \mathbf{A}_{k-1}^T \mathbf{Q}_k^{-1} \mathbf{A}_{k-1} \quad (46a)$$

$$\mathbf{L}_{k-1} \mathbf{d}_{k-1} = \mathbf{q}_{k-1} - \mathbf{A}_{k-1}^T \mathbf{Q}_k^{-1} \mathbf{v}_k \quad (46b)$$

$$\mathbf{L}_{k,k-1} \mathbf{L}_{k-1}^T = -\mathbf{Q}_k^{-1} \mathbf{A}_{k-1} \quad (46c)$$

$$\mathbf{I}_k = -\mathbf{L}_{k,k-1} \mathbf{L}_{k,k-1}^T + \mathbf{Q}_k^{-1} + \mathbf{C}_k^T \mathbf{R}_k^{-1} \mathbf{C}_k \quad (46d)$$

$$\mathbf{q}_k = -\mathbf{L}_{k,k-1} \mathbf{d}_{k-1} + \mathbf{Q}_k^{-1} \mathbf{v}_k + \mathbf{C}_k^T \mathbf{R}_k^{-1} \mathbf{y}_k \quad (46e)$$

- the **backward pass**, $k = K \dots 1$, is

$$\mathbf{L}_{k-1}^T \hat{\mathbf{x}}_{k-1} = -\mathbf{L}_{k,k-1}^T \hat{\mathbf{x}}_k + \mathbf{d}_{k-1} \quad (47)$$

- these are initialized with

$$\mathbf{I}_0 = \check{\mathbf{P}}_0^{-1} + \mathbf{C}_0^T \mathbf{R}_0^{-1} \mathbf{C}_0 \quad (48a)$$

$$\mathbf{q}_0 = \check{\mathbf{P}}_0^{-1} \check{\mathbf{x}}_0 + \mathbf{C}_0^T \mathbf{R}_0^{-1} \mathbf{y}_0 \quad (48b)$$

$$\hat{\mathbf{x}}_K = \mathbf{L}_K^{-T} \mathbf{d}_K \quad (48c)$$

Cholesky solution \mapsto Rauch-Tung-Striebel smoother

- our Cholesky approach will certainly work, and is appealing because it follows easily from the batch setup and the block-tridiagonal sparsity of the LHS
- historically, however, the **Rauch-Tung-Striebel (RTS) smoother** equations constitute the canonical form
- the forwards pass of the RTS smoother is called the Kalman filter
- the Cholesky and RTS equations are algebraically equivalent, but there is quite a lot of algebra needed to show this (see the book)

RTS smoother

- forward pass (a.k.a., Kalman filter), $k = 1 \dots K$:

$$\check{\mathbf{P}}_{k,f} = \mathbf{A}_{k-1} \hat{\mathbf{P}}_{k-1,f} \mathbf{A}_{k-1}^T + \mathbf{Q}_k \quad (49a)$$

$$\check{\mathbf{x}}_{k,f} = \mathbf{A}_{k-1} \hat{\mathbf{x}}_{k-1,f} + \mathbf{v}_k \quad (49b)$$

$$\mathbf{K}_k = \check{\mathbf{P}}_{k,f} \mathbf{C}_k^T (\mathbf{C}_k \check{\mathbf{P}}_{k,f} \mathbf{C}_k^T + \mathbf{R}_k)^{-1} \quad (49c)$$

$$\hat{\mathbf{P}}_{k,f} = (\mathbf{1} - \mathbf{K}_k \mathbf{C}_k) \check{\mathbf{P}}_{k,f} \quad (49d)$$

$$\hat{\mathbf{x}}_{k,f} = \check{\mathbf{x}}_{k,f} + \mathbf{K}_k (\mathbf{y}_k - \mathbf{C}_k \check{\mathbf{x}}_{k,f}) \quad (49e)$$

- backward pass, $k = K \dots 1$:

$$\hat{\mathbf{x}}_{k-1} = \hat{\mathbf{x}}_{k-1,f} + \left(\hat{\mathbf{P}}_{k-1,f} \mathbf{A}_{k-1}^T \check{\mathbf{P}}_{k,f}^{-1} \right) (\hat{\mathbf{x}}_k - \check{\mathbf{x}}_{k,f}) \quad (50a)$$

$$\begin{aligned} \hat{\mathbf{P}}_{k-1} = \hat{\mathbf{P}}_{k-1,f} + \left(\hat{\mathbf{P}}_{k-1,f} \mathbf{A}_{k-1}^T \check{\mathbf{P}}_{k,f}^{-1} \right) & \left(\hat{\mathbf{P}}_k - \check{\mathbf{P}}_{k,f} \right) \\ & \times \left(\hat{\mathbf{P}}_{k-1,f} \mathbf{A}_{k-1}^T \check{\mathbf{P}}_{k,f}^{-1} \right)^T \end{aligned} \quad (50b)$$

- these are initialized with

$$\hat{\mathbf{P}}_{0,f} = (\mathbf{1} - \mathbf{K}_0 \mathbf{C}_0) \check{\mathbf{P}}_0, \quad \hat{\mathbf{x}}_{0,f} = \check{\mathbf{x}}_0 + \mathbf{K}_0 (\mathbf{y}_0 - \mathbf{C}_0 \check{\mathbf{x}}_0) \quad (51a)$$

$$\hat{\mathbf{P}}_K = \hat{\mathbf{P}}_{K,f}, \quad \hat{\mathbf{x}}_K = \hat{\mathbf{x}}_{K,f} \quad (51b)$$

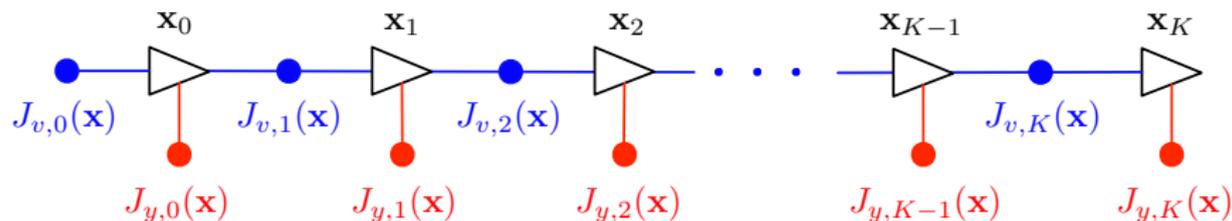
Take-home message

The full Bayesian and MAP approaches result in the **same system of linear equations** because the Bayesian posterior is exactly Gaussian and the mean and mode (i.e., maximum) of a Gaussian are one and the same.

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \hat{\mathbf{x}} = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z} \quad (52)$$

Either the Cholesky or RTS smoother equations can be used to solve this system of equations **exactly** and **efficiently**. These are also equivalent to **Gaussian belief propagation**.

Take-home message



$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) = \arg \min_{\mathbf{x}} \sum_{k=0}^K \left(J_{v,k}(\mathbf{x}) + J_{y,k}(\mathbf{x}) \right)$$

motion prior measurements

$$\underbrace{(\mathbf{A}^{-T} \mathbf{Q}^{-1} \mathbf{A}^{-1} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C})}_{\hat{\mathbf{P}}^{-1}} \hat{\mathbf{x}} = \mathbf{A}^{-T} \mathbf{Q}^{-1} \mathbf{v} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{y}$$

block-tridiagonal

Existence and uniqueness

- how do we know there will be a solution to our batch equations?

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \hat{\mathbf{x}} = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z} \quad (53)$$

- in general, linear systems have 0, 1, or infinitely many solutions
- for a unique solution the LHS might be invertible or equivalently

$$\text{rank}(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) = N(K + 1) \quad (54)$$

because we have $\dim \mathbf{x} = N(K + 1)$

- we can drop \mathbf{W}^{-1} as long as it's invertible – it is, it's positive definite by assumption – the test becomes

$$\text{rank}(\mathbf{H}^T) = N(K + 1) \quad (55)$$

Case (i): we have knowledge of the initial state

- we see that

$$\begin{aligned} & \text{rank } \mathbf{H}^T \\ &= \text{rank} \left[\begin{array}{cccc|cccc} \mathbf{1} & -\mathbf{A}_0^T & & & \mathbf{C}_0^T & & & \\ & \mathbf{1} & -\mathbf{A}_1^T & & & \mathbf{C}_1^T & & \\ & & \mathbf{1} & \ddots & & & \mathbf{C}_2^T & \\ & & & \ddots & -\mathbf{A}_{K-1}^T & & & \ddots \\ & & & & \mathbf{1} & & & \mathbf{C}_K^T \end{array} \right] \end{aligned} \quad (56)$$

- is exactly in **block-row-echelon form**, meaning there are $N(K + 1)$ 'leading ones'
- this means all the rows are linearly independent so the matrix is full rank and our **solution is unique**
- intuitively, the prior already provides a complete solution and the measurements nudge it

Case (ii): no knowledge of the initial state

- it might be that we have no idea about the initial state of the system
- this removes the first column in \mathbf{H}^T and our condition becomes

$$\begin{aligned} & \text{rank } \mathbf{H}^T \\ = & \text{rank} \left[\begin{array}{cccc|cccc} -\mathbf{A}_0^T & & & & \mathbf{C}_0^T & & & & \\ \mathbf{1} & -\mathbf{A}_1^T & & & & \mathbf{C}_1^T & & & \\ & \mathbf{1} & \ddots & & & & \mathbf{C}_2^T & & \\ & & \ddots & -\mathbf{A}_{K-1}^T & & & & \ddots & \\ & & & \mathbf{1} & & & & & \mathbf{C}_K^T \end{array} \right] \quad (57) \end{aligned}$$

- it's no longer obvious that the solution is unique

Case (ii): no knowledge of the initial state

- moving the top block-row to the bottom does not change the rank

$$\begin{aligned} & \text{rank } \mathbf{H}^T \\ &= \text{rank} \left[\begin{array}{ccc|ccc} \mathbf{1} & -\mathbf{A}_1^T & & & \mathbf{C}_1^T & \\ & \mathbf{1} & \ddots & & & \mathbf{C}_2^T \\ & & \ddots & -\mathbf{A}_{K-1}^T & & \ddots \\ & & & \mathbf{1} & & \mathbf{C}_K^T \\ \hline -\mathbf{A}_0^T & & & & \mathbf{C}_0^T & \end{array} \right] \quad (58) \end{aligned}$$

- again without changing the rank, we can add to the bottom block-row, \mathbf{A}_0^T times the first block-row, $\mathbf{A}_0^T \mathbf{A}_1^T$ times the second block-row, \dots , and $\mathbf{A}_0^T \dots \mathbf{A}_{K-1}^T$ times the K th block-row

Case (ii): no knowledge of the initial state

- the result of this last step is

$$\begin{aligned} & \text{rank } \mathbf{H}^T \\ = & \text{rank} \left[\begin{array}{cccc|cccc} 1 & -\mathbf{A}_1^T & & & & & & \mathbf{C}_1^T \\ & 1 & \ddots & & & & & \mathbf{C}_1^T \\ & & \ddots & & & & & \vdots \\ & & & -\mathbf{A}_{K-1}^T & & & & \mathbf{C}_K^T \\ & & & 1 & & & & \mathbf{C}_K^T \\ \hline & & & & \mathbf{C}_0^T & \mathbf{A}_0^T \mathbf{C}_1^T & \mathbf{A}_0^T \mathbf{A}_1^T \mathbf{C}_2^T & \dots & \mathbf{A}_0^T \dots \mathbf{A}_{K-1}^T \mathbf{C}_K^T \end{array} \right] \end{aligned} \quad (59)$$

- the upper-left block is full rank, NK , since every row has a ‘leading one’
- all that remains is to determine if the bottom-right block is full rank:

$$\text{rank} \left[\mathbf{C}_0^T \quad \mathbf{A}_0^T \mathbf{C}_1^T \quad \mathbf{A}_0^T \mathbf{A}_1^T \mathbf{C}_2^T \quad \dots \quad \mathbf{A}_0^T \dots \mathbf{A}_{K-1}^T \mathbf{C}_K^T \right] = N \quad (60)$$

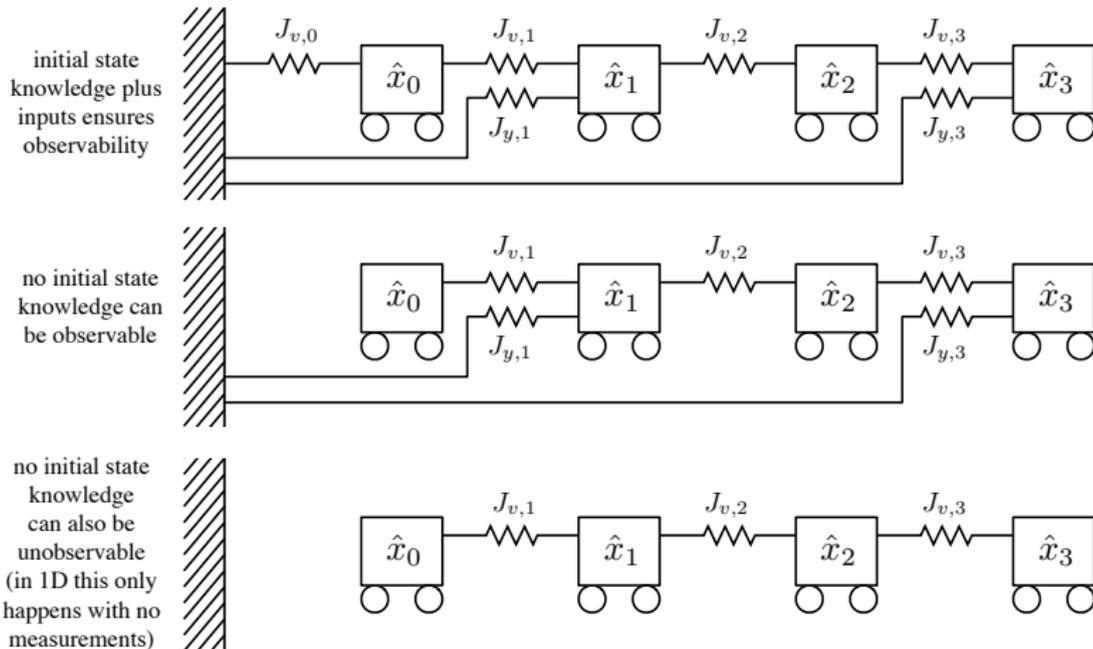
Case (ii): no knowledge of the initial state – time-invariant

- if we further assume the system is time-invariant such that for all k we have $\mathbf{A}_k = \mathbf{A}$ and $\mathbf{C}_k = \mathbf{C}$ and we make the not-too-restrictive assumption that $K \gg N$, we may further simplify this condition
- the **Cayley-Hamilton theorem** says square matrices satisfy their own characteristic equation, so powers of \mathbf{A} greater than or equal to N can be written as linear combinations of the lower powers
- this let's us keep only the first N block-rows, reducing the test to

$$\text{rank} \begin{bmatrix} \mathbf{C} \\ \mathbf{C}\mathbf{A} \\ \vdots \\ \mathbf{C}\mathbf{A}^{(N-1)} \end{bmatrix} = N \quad (61)$$

- this is precisely the test for **observability** – that we can reconstruct the initial state from a finite number of measurements

Mass-spring analogy



Knowing when we don't know

- we said at the beginning of the lecture that we'd like to bookkeep all of the uncertainties in our estimate, but the linear system just tells us the mean, right?
- wrong, we can interpret it in the following way:

$$\underbrace{(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})}_{\substack{\text{inverse} \\ \text{covariance}}} \underbrace{\hat{\mathbf{x}}}_{\text{mean}} = \underbrace{\mathbf{H}^T \mathbf{W}^{-1} \mathbf{z}}_{\substack{\text{information} \\ \text{vector}}} \quad (62)$$

- so the covariance of our estimate is

$$\hat{\mathbf{P}} = (\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \quad (63)$$

- our (Bayesian) estimate is

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\hat{\mathbf{x}}, \hat{\mathbf{P}}) \quad (64)$$

Knowing how to know when we don't know

- to see this, we can directly take the expectation of the estimate:

$$\mathbf{x} - \underbrace{(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{W}^{-1} \mathbf{z}}_{E[\mathbf{x}]} = (\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{W}^{-1} \underbrace{(\mathbf{H} \mathbf{x} - \mathbf{z})}_{\mathbf{s}} \quad (65)$$

where $\mathbf{s} = \begin{bmatrix} \mathbf{w} \\ \mathbf{n} \end{bmatrix}$

- in this case we have

$$\hat{\mathbf{P}} = E \left[(\mathbf{x} - E[\mathbf{x}]) (\mathbf{x} - E[\mathbf{x}])^T \right] \quad (66a)$$

$$= (\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{W}^{-1} \underbrace{E[\mathbf{s} \mathbf{s}^T]}_{\mathbf{W}} \mathbf{W}^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \quad (66b)$$

$$= (\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H})^{-1} \quad (66c)$$

Lecture Summary

- we have looked at how to estimate a whole trajectory using initial state knowledge, a sequence of inputs, and a sequence of measurements for linear systems corrupted by Gaussian noise
- our Bayesian method involves solving a linear system of equations for the mean and also builds the covariance
- we know how to solve this linear system very efficiently by **exploiting the sparsity** of the matrices
- we have a test for when the solution is unique
- the batch approach can only be used offline since it is **acausal** (it uses future data to estimate past states)
- the forward pass of the RTS smoother is the recursive **Kalman filter**, which can be used online
- **Lecture 2: Nonlinearities and Lie Groups**

Lecture 2: Nonlinearities and Lie Groups

A Short Course in State Estimation

Timothy D. Barfoot

University of Toronto

Copyright © 2024

Lecture Outline

Lecture 1: Linear-Gaussian Estimation

Lecture 2: Nonlinearities and Lie Groups

- Problem Setup

- Gauss-Newton

- Sparsity

- Lie Groups and Algebras

- Perturbations and Optimization

- Example: Point-Cloud Alignment

Lecture 3: Continuous-Time Estimation

Lecture 4: Beyond MAP – Variational Inference

System

- we define our system using the following **nonlinear** models:

$$\text{motion model: } \mathbf{x}_k = \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{v}_k, \mathbf{w}_k), \quad k = 1 \dots K \quad (67a)$$

$$\text{observation model: } \mathbf{y}_k = \mathbf{g}(\mathbf{x}_k, \mathbf{n}_k), \quad k = 0 \dots K \quad (67b)$$

where k is again the discrete-time index and K its maximum

- the variables have the following meanings:

$$\text{system state : } \mathbf{x}_k \in \mathbb{R}^N$$

$$\text{initial state : } \mathbf{x}_0 \in \mathbb{R}^N \sim \mathcal{N}(\check{\mathbf{x}}_0, \check{\mathbf{P}}_0)$$

$$\text{input : } \mathbf{v}_k \in \mathbb{R}^N$$

$$\text{process noise : } \mathbf{w}_k \in \mathbb{R}^N \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k)$$

$$\text{measurement : } \mathbf{y}_k \in \mathbb{R}^M$$

$$\text{measurement noise : } \mathbf{n}_k \in \mathbb{R}^M \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$$

Nonlinear errors

- we aim to use the **MAP** approach to state estimation, but now for **nonlinear** systems
- we define the **errors** with respect to the prior and measurements to be

$$\mathbf{e}_{v,k}(\mathbf{x}) = \begin{cases} \check{\mathbf{x}}_0 - \mathbf{x}_0, & k = 0 \\ \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{v}_k, \mathbf{0}) - \mathbf{x}_k, & k = 1 \dots K \end{cases} \quad (68a)$$

$$\mathbf{e}_{y,k}(\mathbf{x}) = \mathbf{y}_k - \mathbf{g}(\mathbf{x}_k, \mathbf{0}), \quad k = 0 \dots K \quad (68b)$$

- as usual, there is one error term for every measurement, every input, and one for the initial state knowledge

MAP cost terms

- the contributions to the objective function are

$$J_{v,k}(\mathbf{x}) = \frac{1}{2} \mathbf{e}_{v,k}(\mathbf{x})^T \mathbf{W}_{v,k}^{-1} \mathbf{e}_{v,k}(\mathbf{x}) \quad (69a)$$

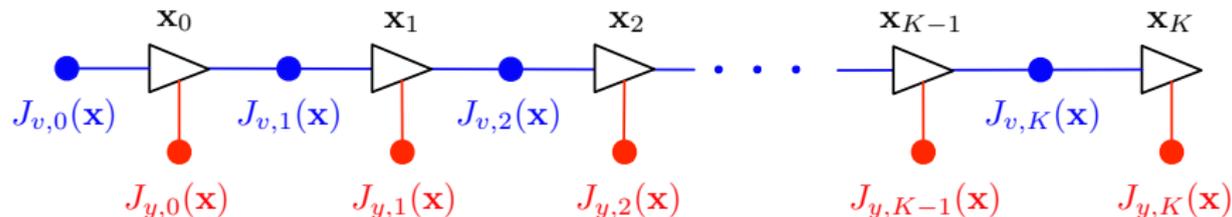
$$J_{y,k}(\mathbf{x}) = \frac{1}{2} \mathbf{e}_{y,k}(\mathbf{x})^T \mathbf{W}_{y,k}^{-1} \mathbf{e}_{y,k}(\mathbf{x}) \quad (69b)$$

- the overall objective function is then

$$J(\mathbf{x}) = \sum_{k=0}^K (J_{v,k}(\mathbf{x}) + J_{y,k}(\mathbf{x})) \quad (70)$$

- by choosing $\mathbf{W}_{v,k}$ and $\mathbf{W}_{y,k}$ to be the inverse covariances of the measurement noises, minimizing the objective function is equivalent to maximizing the joint likelihood of all the data

MAP as a factor graph



$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) = \arg \min_{\mathbf{x}} \sum_{k=0}^K \left(J_{v,k}(\mathbf{x}) + J_{y,k}(\mathbf{x}) \right)$$

motion prior measurements

Stacking things up

– we further define

$$\mathbf{e}(\mathbf{x}) = \begin{bmatrix} \mathbf{e}_v(\mathbf{x}) \\ \mathbf{e}_y(\mathbf{x}) \end{bmatrix}, \quad \mathbf{e}_v(\mathbf{x}) = \begin{bmatrix} \mathbf{e}_{v,0}(\mathbf{x}) \\ \vdots \\ \mathbf{e}_{v,K}(\mathbf{x}) \end{bmatrix}, \quad \mathbf{e}_y(\mathbf{x}) = \begin{bmatrix} \mathbf{e}_{y,0}(\mathbf{x}) \\ \vdots \\ \mathbf{e}_{y,K}(\mathbf{x}) \end{bmatrix} \quad (71a)$$

$$\mathbf{W} = \text{diag}(\mathbf{W}_v, \mathbf{W}_y), \quad \mathbf{W}_v = \text{diag}(\mathbf{W}_{v,0}, \dots, \mathbf{W}_{v,K}), \quad (71b)$$

$$\mathbf{W}_y = \text{diag}(\mathbf{W}_{y,0}, \dots, \mathbf{W}_{y,K}) \quad (71c)$$

so that the objective function can be written as

$$J(\mathbf{x}) = \frac{1}{2} \mathbf{e}(\mathbf{x})^T \mathbf{W}^{-1} \mathbf{e}(\mathbf{x}) \quad (72)$$

Standard nonlinear squared error form

- we can further define the modified error term,

$$\mathbf{u}(\mathbf{x}) = \mathbf{L} \mathbf{e}(\mathbf{x}) \quad (73)$$

where $\mathbf{L}^T \mathbf{L} = \mathbf{W}^{-1}$ (i.e., from a **Cholesky decomposition** since \mathbf{W} is symmetric positive-definite)

- using these definitions, we can write the objective function simply as

$$J(\mathbf{x}) = \frac{1}{2} \mathbf{u}(\mathbf{x})^T \mathbf{u}(\mathbf{x}) \quad (74)$$

- this is precisely in a quadratic form, but not with respect to the design variables, \mathbf{x}
- due to the nonlinearities, this is typically a **nonconvex** function, which can have **local minima**

Gauss-Newton

- the **Gauss-Newton** optimization method iteratively approximates the objective function as a quadratic, which can be analytically minimized at each iteration – result is a **local** minimization of $J(\mathbf{x})$
- a ‘short-cut’ to the Gauss-Newton method is to start with a first-order Taylor expansion of $\mathbf{u}(\mathbf{x})$:

$$\mathbf{u}(\mathbf{x}_{\text{op}} + \delta\mathbf{x}) \approx \mathbf{u}(\mathbf{x}_{\text{op}}) + \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) \delta\mathbf{x} \quad (75)$$

- substituting into J we have

$$\begin{aligned} & J(\mathbf{x}_{\text{op}} + \delta\mathbf{x}) \\ & \approx \frac{1}{2} \left(\mathbf{u}(\mathbf{x}_{\text{op}}) + \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) \delta\mathbf{x} \right)^T \left(\mathbf{u}(\mathbf{x}_{\text{op}}) + \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) \delta\mathbf{x} \right) \end{aligned} \quad (76)$$

Gauss-Newton

- minimizing with respect to $\delta \mathbf{x}$ gives

$$\begin{aligned} \frac{\partial J(\mathbf{x}_{\text{op}} + \delta \mathbf{x})}{\partial \delta \mathbf{x}} &= \left(\mathbf{u}(\mathbf{x}_{\text{op}}) + \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) \delta \mathbf{x}^* \right)^T \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) = \mathbf{0} \\ \Rightarrow \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right)^T \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right) \delta \mathbf{x}^* &= - \left(\left. \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \right)^T \mathbf{u}(\mathbf{x}_{\text{op}}) \quad (77) \end{aligned}$$

- we apply the update to our initial guess using

$$\mathbf{x}_{\text{op}} \leftarrow \mathbf{x}_{\text{op}} + \delta \mathbf{x}^* \quad (78)$$

and iterate to convergence

- modifications such as **line search**, **Levenberg-Marquardt**, and **trust-region methods** can be used to improve convergence

Maximum a posteriori

- getting back to our specific setup, we recall that

$$\mathbf{u}(\mathbf{x}) = \mathbf{L} \mathbf{e}(\mathbf{x}) \quad (79)$$

with \mathbf{L} a constant

- we substitute this into the Gauss-Newton update to see that in terms of the error, $\mathbf{e}(\mathbf{x})$, we have

$$(\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}) \delta \mathbf{x}^* = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{e}(\mathbf{x}_{\text{op}}) \quad (80)$$

with

$$\mathbf{H} = - \left. \frac{\partial \mathbf{e}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_{\text{op}}} \quad (81)$$

and where we have used $\mathbf{L}^T \mathbf{L} = \mathbf{W}^{-1}$

MAP

– at the individual error level, the linearized approximations are:

$$\mathbf{e}_{v,k}(\mathbf{x}_{\text{op}} + \delta\mathbf{x}) \approx \begin{cases} \mathbf{e}_{v,0}(\mathbf{x}_{\text{op}}) - \delta\mathbf{x}_0, & k = 0 \\ \mathbf{e}_{v,k}(\mathbf{x}_{\text{op}}) + \mathbf{F}_{k-1}\delta\mathbf{x}_{k-1} - \delta\mathbf{x}_k, & k = 1 \dots K \end{cases} \quad (82)$$

$$\mathbf{e}_{y,k}(\mathbf{x}_{\text{op}} + \delta\mathbf{x}) \approx \mathbf{e}_{y,k}(\mathbf{x}_{\text{op}}) - \mathbf{G}_k\delta\mathbf{x}_k, \quad k = 0 \dots K \quad (83)$$

where

$$\mathbf{e}_{v,k}(\mathbf{x}_{\text{op}}) \approx \begin{cases} \check{\mathbf{x}}_0 - \mathbf{x}_{\text{op},0}, & k = 0 \\ \mathbf{f}(\mathbf{x}_{\text{op},k-1}, \mathbf{v}_k, \mathbf{0}) - \mathbf{x}_{\text{op},k}, & k = 1 \dots K \end{cases} \quad (84)$$

$$\mathbf{e}_{y,k}(\mathbf{x}_{\text{op}}) \approx \mathbf{y}_k - \mathbf{g}(\mathbf{x}_{\text{op},k}, \mathbf{0}), \quad k = 0 \dots K \quad (85)$$

$$\mathbf{F}_{k-1} = \left. \frac{\partial \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{v}_k, \mathbf{w}_k)}{\partial \mathbf{x}_{k-1}} \right|_{\mathbf{x}_{\text{op},k-1}, \mathbf{v}_k, \mathbf{0}}, \quad \mathbf{G}_k = \left. \frac{\partial \mathbf{g}(\mathbf{x}_k, \mathbf{n}_k)}{\partial \mathbf{x}_k} \right|_{\mathbf{x}_{\text{op},k}, \mathbf{0}} \quad (86)$$

Exploiting sparsity

- just as in the linear-Gaussian case, we must solve a large linear system of equations where the left-hand side is **block-tridiagonal** – only now we must do this **at each iteration** of Gauss-Newton
- we can modify the **Cholesky** and **RTS smoothers** to help us solve this linear system of equations
 - if we use the RTS smoother in a single forward and backward pass where the linearization point is improved at each timestep, this is referred to as the **extended RTS smoother** – the forward pass of this is the **extended Kalman filter**
 - if iterate over the whole trajectory, the RTS smoother can efficiently implement the full Gauss-Newton method
 - there are schemes in-between these two extremes such as the **iterated EKF** and **sliding-window filters**

Iterating is key

Gauss-Newton iterates over the entire trajectory, but runs offline and not in constant time

\mathbf{x}_0 \mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \cdots \mathbf{x}_{k-2} \mathbf{x}_{k-1} \mathbf{x}_k \mathbf{x}_{k+1} \mathbf{x}_{k+2} \cdots \mathbf{x}_K



Sliding-window filters iterate over several timesteps at once, run online and in constant time

\mathbf{x}_0 \mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \cdots \mathbf{x}_{k-2} \mathbf{x}_{k-1} \mathbf{x}_k \mathbf{x}_{k+1} \mathbf{x}_{k+2} \cdots \mathbf{x}_K



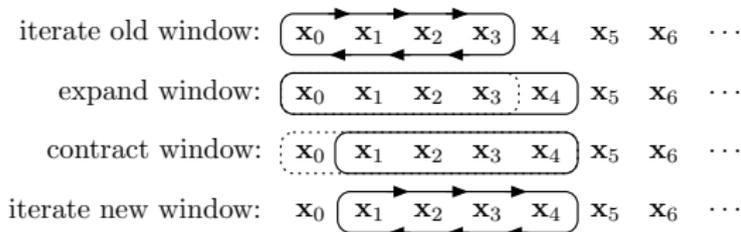
IEKF iterates at only one timestep at a time, but runs online and in constant time

\mathbf{x}_0 \mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \cdots \mathbf{x}_{k-2} \mathbf{x}_{k-1} \mathbf{x}_k \mathbf{x}_{k+1} \mathbf{x}_{k+2} \cdots \mathbf{x}_K



Sliding-window filters

- **sliding-window filters** are a very common way to strike a nice balance (accuracy vs. computational cost) between full batch methods and basic filters:



- while there are several variants, most SWFs use **marginalization** in the window-contraction step:

$$\underbrace{\begin{bmatrix} \mathbf{A}_{00} & \mathbf{A}_{10}^T \\ \mathbf{A}_{10} & \mathbf{A}_{11} & \mathbf{A}_{21}^T \\ & \mathbf{A}_{21} & \mathbf{A}_{22} & \mathbf{A}_{32}^T \\ & & \mathbf{A}_{32} & \mathbf{A}_{33} & \mathbf{A}_{43}^T \\ & & & \mathbf{A}_{43} & \mathbf{A}_{44} \end{bmatrix}}_{\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}} \underbrace{\begin{bmatrix} \delta x_0^* \\ \delta x_1^* \\ \delta x_2^* \\ \delta x_3^* \\ \delta x_4^* \end{bmatrix}}_{\delta x^*} = \underbrace{\begin{bmatrix} \mathbf{b}_0 \\ \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{b}_3 \\ \mathbf{b}_4 \end{bmatrix}}_{\mathbf{H}^T \mathbf{W}^{-1} \mathbf{e}} \Rightarrow \underbrace{\begin{bmatrix} \mathbf{A}_{11} - \mathbf{A}_{10} \mathbf{A}_{00}^{-1} \mathbf{A}_{10}^T & \mathbf{A}_{21}^T \\ & \mathbf{A}_{22} & \mathbf{A}_{32}^T \\ & \mathbf{A}_{32} & \mathbf{A}_{33} & \mathbf{A}_{43}^T \\ & & \mathbf{A}_{43} & \mathbf{A}_{44} \end{bmatrix}}_{\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H}} \underbrace{\begin{bmatrix} \delta x_1^* \\ \delta x_2^* \\ \delta x_3^* \\ \delta x_4^* \end{bmatrix}}_{\delta x^*} = \underbrace{\begin{bmatrix} \mathbf{b}_1 - \mathbf{A}_{10} \mathbf{A}_{00}^{-1} \mathbf{b}_0 \\ \mathbf{b}_2 \\ \mathbf{b}_3 \\ \mathbf{b}_4 \end{bmatrix}}_{\mathbf{H}^T \mathbf{W}^{-1} \mathbf{e}}$$

We still need more tools

- consider the problem of **aligning two point-clouds**, \mathbf{y}_j and \mathbf{p}_j
- we define our error term for each point as

$$\mathbf{e}_j = \mathbf{y}_j - \mathbf{T}\mathbf{p}_j \quad (90)$$

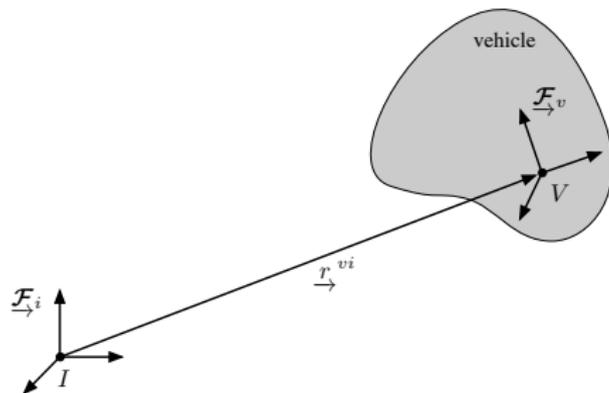
and our objective function (with no motion prior) as

$$J(\mathbf{T}) = \frac{1}{2} \sum_{j=1}^M w_j \mathbf{e}_j^T \mathbf{e}_j = \frac{1}{2} \sum_{j=1}^M w_j (\mathbf{y}_j - \mathbf{T}\mathbf{p}_j)^T (\mathbf{y}_j - \mathbf{T}\mathbf{p}_j) \quad (91)$$

where $w_j > 0$ are the usual scalar weights

- our state in this problem, \mathbf{T} , is the **unknown pose** between the two point-clouds
- \mathbf{T} is not a typical vector quantity – it has **constraints** on its form that we need to account for in the optimization problem

Pose



- the **pose** of one reference frame with respect to another has six degrees of freedom:
 - three in **translation**
 - three in **rotation**
- if we know the pose, $\{\mathbf{r}_i^{vi}, \mathbf{C}_{iv}\}$, we can transform the coordinates of a point, P , from one frame to another:

$$\mathbf{r}_i^{pi} = \mathbf{C}_{iv}\mathbf{r}_v^{pv} + \mathbf{r}_i^{vi} \quad (92)$$

Transformation matrices

- we can combine the translation and rotation of a pose into a convenient form called the (4×4) **transformation matrix**:

$$\begin{bmatrix} \mathbf{r}_i^{pi} \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{C}_{iv} & \mathbf{r}_i^{vi} \\ \mathbf{0}^T & 1 \end{bmatrix}}_{\mathbf{T}_{iv}} \begin{bmatrix} \mathbf{r}_v^{pv} \\ 1 \end{bmatrix} \quad (93)$$

- we see that this allows us to easily transform points from one frame to another in so-called (4×1) **homogenous** point representation:

$$\begin{bmatrix} \mathbf{r} \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \quad (94)$$

which just has an **extra 1** at the bottom

Transformation matrices

- transformation matrices are 4×4 and always have this special structure:

$$\mathbf{T} = \begin{bmatrix} \mathbf{C} & \mathbf{r} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (95)$$

- they have 16 parameters but only 6 degrees of freedom and therefore must have 10 constraints
- 6 constraints come from $\mathbf{C}^T \mathbf{C} = \mathbf{1}$ and the other 4 come from the fact that the bottom row is always $(0, 0, 0, 1)$
- we can compound transformation matrices (just like rotation matrices):

$$\mathbf{T}_{iv} = \mathbf{T}_{ia} \mathbf{T}_{ab} \mathbf{T}_{bv} \quad (96)$$

and the structure always holds (more on this later); **order matters**

Matrix Lie groups

- **rotation** and **transformation** matrices show up in our kinematic (i.e., motion) and sensor (i.e., observation) models
- neither behaves like a vector, yet all of our estimation tools assume the state is a vector
- it turns out that the sets of rotations and poses are not vector spaces, but another type of mathematical object called **matrix Lie groups**
- we will use the rest of this lecture to learn about this in the hope that it will guide us in the estimation of rotations and poses
- **spoiler**: we will use properties of the Lie-group structure to maintain our estimation algorithms as **unconstrained** optimizations

Special orthogonal group

- the set of rotations is called the **special orthogonal group**:

$$SO(3) = \{ \mathbf{C} \in \mathbb{R}^{3 \times 3} \mid \mathbf{C}\mathbf{C}^T = \mathbf{1}, \det \mathbf{C} = 1 \} \quad (97)$$

- the $\mathbf{C}\mathbf{C}^T = \mathbf{1}$ orthogonality condition is needed to impose 6 constraints on the 9-parameter rotation matrix, reducing the degrees of freedom to 3
- noticing that

$$(\det \mathbf{C})^2 = \det (\mathbf{C}\mathbf{C}^T) = \det \mathbf{1} = 1 \quad (98)$$

we have that $\det \mathbf{C} = \pm 1$, allowing for two possibilities

- choosing $\det \mathbf{C} = 1$ ensures that we have a **proper rotation**, and it's this additional property that makes the $O(3)$ orthogonality group **special**

Special Euclidean group

- the set of transformation matrices representing poses is called the **special Euclidean group**:

$$SE(3) = \left\{ \mathbf{T} = \begin{bmatrix} \mathbf{C} & \mathbf{r} \\ \mathbf{0}^T & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \mid \mathbf{C} \in SO(3), \mathbf{r} \in \mathbb{R}^3 \right\} \quad (99)$$

Matrix Lie groups

- both $SO(3)$ and $SE(3)$ are **matrix Lie groups**
- to be a **group** they must have an operator to combine elements that satisfies 4 properties: closure, associativity, identity, invertibility
- to be a **Lie** group the operator must be 'smooth'
- to be a **matrix** Lie group the elements must be matrices and the operator matrix multiplication

property	$SO(3)$	$SE(3)$
closure	$\mathbf{C}_1, \mathbf{C}_2 \in SO(3)$ $\Rightarrow \mathbf{C}_1 \mathbf{C}_2 \in SO(3)$	$\mathbf{T}_1, \mathbf{T}_2 \in SE(3)$ $\Rightarrow \mathbf{T}_1 \mathbf{T}_2 \in SE(3)$
associativity	$\mathbf{C}_1 (\mathbf{C}_2 \mathbf{C}_3) = (\mathbf{C}_1 \mathbf{C}_2) \mathbf{C}_3$ $= \mathbf{C}_1 \mathbf{C}_2 \mathbf{C}_3$	$\mathbf{T}_1 (\mathbf{T}_2 \mathbf{T}_3) = (\mathbf{T}_1 \mathbf{T}_2) \mathbf{T}_3$ $= \mathbf{T}_1 \mathbf{T}_2 \mathbf{T}_3$
identity	$\mathbf{C}, \mathbf{1} \in SO(3)$ $\Rightarrow \mathbf{C} \mathbf{1} = \mathbf{1} \mathbf{C} = \mathbf{C}$	$\mathbf{T}, \mathbf{1} \in SE(3)$ $\Rightarrow \mathbf{T} \mathbf{1} = \mathbf{1} \mathbf{T} = \mathbf{T}$
invertibility	$\mathbf{C} \in SO(3)$ $\Rightarrow \mathbf{C}^{-1} \in SO(3)$	$\mathbf{T} \in SE(3)$ $\Rightarrow \mathbf{T}^{-1} \in SE(3)$

Lie algebras

- to every matrix Lie group there is associated a **Lie algebra**, which consists of a vector space, \mathbb{V} , over some field, \mathbb{F} , together with a binary operation, $[\cdot, \cdot]$, called the **Lie bracket** (of the algebra) that satisfies four properties:

closure: $[\mathbf{X}, \mathbf{Y}] \in \mathbb{V}$

bilinearity: $[a\mathbf{X} + b\mathbf{Y}, \mathbf{Z}] = a[\mathbf{X}, \mathbf{Z}] + b[\mathbf{Y}, \mathbf{Z}]$,
 $[\mathbf{Z}, a\mathbf{X} + b\mathbf{Y}] = a[\mathbf{Z}, \mathbf{X}] + b[\mathbf{Z}, \mathbf{Y}]$

alternating: $[\mathbf{X}, \mathbf{X}] = \mathbf{0}$

Jacobi identity: $[\mathbf{X}, [\mathbf{Y}, \mathbf{Z}]] + [\mathbf{Z}, [\mathbf{Y}, \mathbf{X}]] + [\mathbf{Y}, [\mathbf{Z}, \mathbf{X}]] = \mathbf{0}$

for all $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \in \mathbb{V}$ and $a, b \in \mathbb{F}$

Lie algebra: rotations

- the **Lie algebra** associated with $SO(3)$ is given by

$$\text{vector space: } \mathfrak{so}(3) = \{ \Phi = \phi^\wedge \in \mathbb{R}^{3 \times 3} \mid \phi \in \mathbb{R}^3, \}$$

$$\text{field: } \mathbb{R}$$

$$\text{Lie bracket: } [\Phi_1, \Phi_2] = \Phi_1 \Phi_2 - \Phi_2 \Phi_1$$

where

$$\phi^\wedge = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \end{bmatrix}^\wedge = \begin{bmatrix} 0 & -\phi_3 & \phi_2 \\ \phi_3 & 0 & -\phi_1 \\ -\phi_2 & \phi_1 & 0 \end{bmatrix} \in \mathbb{R}^{3 \times 3}, \quad \phi \in \mathbb{R}^3 \quad (100)$$

is the (linear) **skew-symmetric operator**

Lie algebra: poses

- the **Lie algebra** associated with $SE(3)$ is given by

$$\text{vector space: } \mathfrak{se}(3) = \{\Xi = \xi^\wedge \in \mathbb{R}^{4 \times 4} \mid \xi \in \mathbb{R}^6\}$$

$$\text{field: } \mathbb{R}$$

$$\text{Lie bracket: } [\Xi_1, \Xi_2] = \Xi_1 \Xi_2 - \Xi_2 \Xi_1$$

where

$$\xi^\wedge = \begin{bmatrix} \rho \\ \phi \end{bmatrix}^\wedge = \begin{bmatrix} \phi^\wedge & \rho \\ \mathbf{0}^T & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4}, \quad \rho, \phi \in \mathbb{R}^3 \quad (101)$$

- this is an **overloading** of the $(\cdot)^\wedge$ operator from before to take elements of \mathbb{R}^6 and turn them into elements of $\mathbb{R}^{4 \times 4}$; it is still linear

This is getting exponentially more complicated

- ok, so the sets of rotation and transformation matrices are **matrix Lie groups**: $SO(3)$ and $SE(3)$
- each one has an associated **Lie algebra**: $\mathfrak{so}(3)$ and $\mathfrak{se}(3)$
- so what, where is this all going?!
- to get to the next level of understanding, we need a connection between the Lie group and Lie algebra
- that connection is the **exponential map** given by

$$\exp(\mathbf{A}) = \mathbf{1} + \mathbf{A} + \frac{1}{2!}\mathbf{A}^2 + \frac{1}{3!}\mathbf{A}^3 + \dots = \sum_{n=0}^{\infty} \frac{1}{n!}\mathbf{A}^n \quad (102)$$

where $\mathbf{A} \in \mathbb{R}^{M \times M}$ is a square matrix

Exponential rotations

- for **rotations**, we can relate elements of $SO(3)$ to elements of $\mathfrak{so}(3)$ through the exponential map:

$$\mathbf{C} = \exp(\phi^\wedge) = \sum_{n=0}^{\infty} \frac{1}{n!} (\phi^\wedge)^n \quad (103)$$

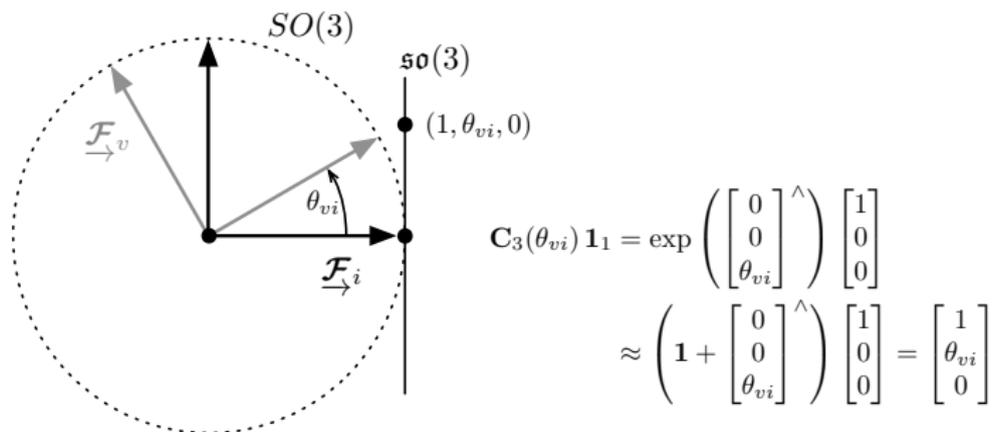
where $\mathbf{C} \in SO(3)$ and $\phi \in \mathbb{R}^3$ (and hence $\phi^\wedge \in \mathfrak{so}(3)$)

- we can also go in the other direction (but not uniquely) using

$$\phi = \ln(\mathbf{C})^\vee \quad (104)$$

- the mapping is **surjective** (or onto), meaning every element of $SO(3)$ can be generated by at least one element of $\mathfrak{so}(3)$
- the non-unique inverse mapping is due to **singularities** – in this case $\phi + 2\pi m$ with m any integer produces the same \mathbf{C}

Tangent space



- the vector space of a Lie algebra is the **tangent space** of the associated Lie group at the identity element of the group, and it completely captures the local structure of the group

Exponential poses

- for **poses**, we can relate elements of $SE(3)$ to elements of $\mathfrak{se}(3)$, again through the exponential map:

$$\mathbf{T} = \exp(\boldsymbol{\xi}^\wedge) = \sum_{n=0}^{\infty} \frac{1}{n!} (\boldsymbol{\xi}^\wedge)^n \quad (105)$$

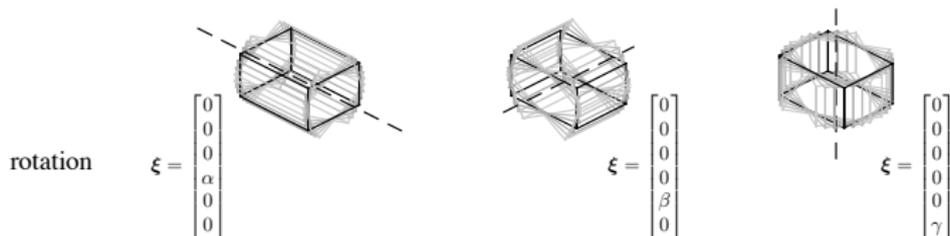
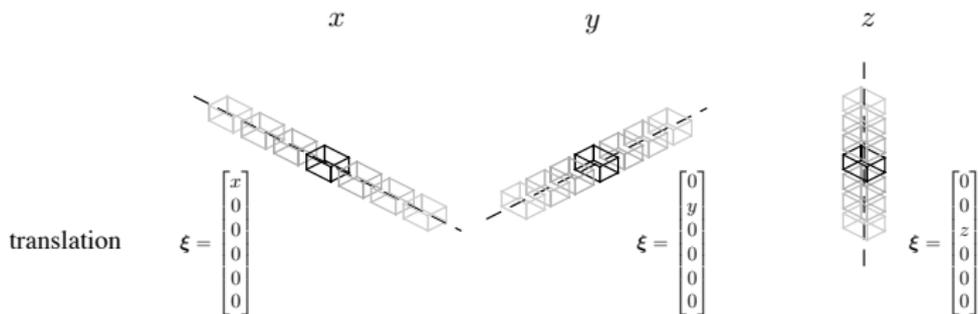
where $\mathbf{T} \in SE(3)$ and $\boldsymbol{\xi} \in \mathbb{R}^6$ (and hence $\boldsymbol{\xi}^\wedge \in \mathfrak{se}(3)$)

- we can also go in the other direction (again, not uniquely) using

$$\boldsymbol{\xi} = \ln(\mathbf{T})^\vee \quad (106)$$

- the exponential map from $\mathfrak{se}(3)$ to $SE(3)$ is also surjective: every $\boldsymbol{\xi} \in \mathbb{R}^6$ maps to some $\mathbf{T} \in SE(3)$ (many-to-one) and every $\mathbf{T} \in SE(3)$ can be generated by at least one $\boldsymbol{\xi} \in \mathbb{R}^6$

Pose change



- varying each component of ξ then using $\mathbf{T} = \exp(\xi^\wedge)$ to transform the points comprising the corners of a rectangular prism

Perturbations

- we now introduce the idea of **perturbations**
- for **vectors**, we usually perturb like this:

$$\mathbf{x} = \underbrace{\bar{\mathbf{x}}}_{\text{'big'}} + \underbrace{\delta\mathbf{x}}_{\text{'small'}} \quad (107)$$

but actually this is an arbitrary choice

- in an **optimization** setting, perturbations are used like this:

$$\mathbf{x} = \underbrace{\bar{\mathbf{x}}}_{\text{initial guess}} + \underbrace{\delta\mathbf{x}}_{\text{optimal update}} \quad (108)$$

- in a **probability** setting, perturbations are used like this:

$$\mathbf{x} = \underbrace{\bar{\mathbf{x}}}_{\text{deterministic}} + \underbrace{\delta\mathbf{x}}_{\text{random noise}} \quad (109)$$

Rotation perturbations

- for **rotations**, we will perturb like this:

$$\mathbf{C} = \underbrace{\delta\mathbf{C}}_{\text{'small'}} \underbrace{\bar{\mathbf{C}}}_{\text{'big'}} \quad (110)$$

- we pick the following **perturbation**

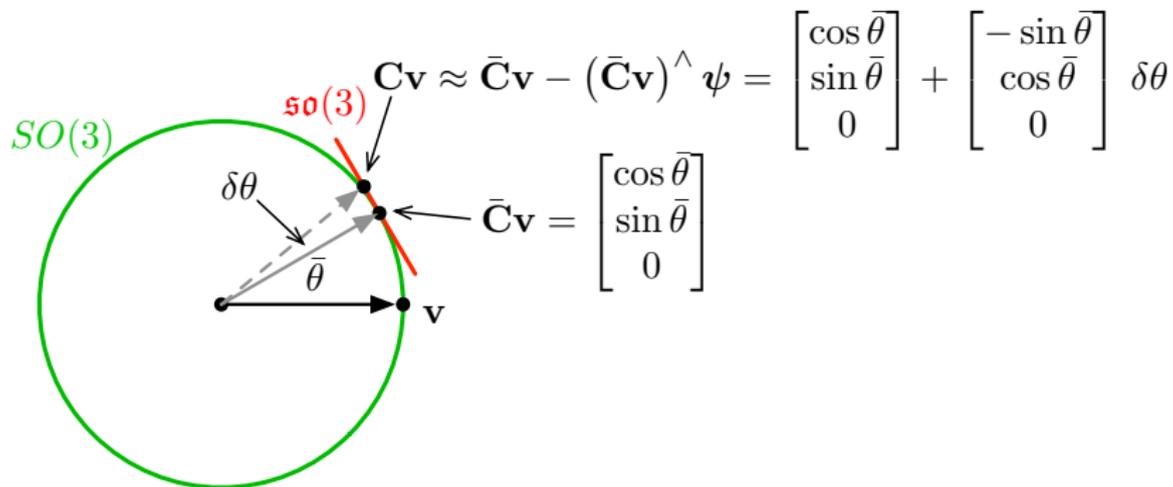
$$\delta\mathbf{C} = \exp(\psi^\wedge) \quad (111)$$

which ensures \mathbf{C} is still a valid rotation (by closure)

- this lets us **linearize** the product of a rotation and point, \mathbf{v} :

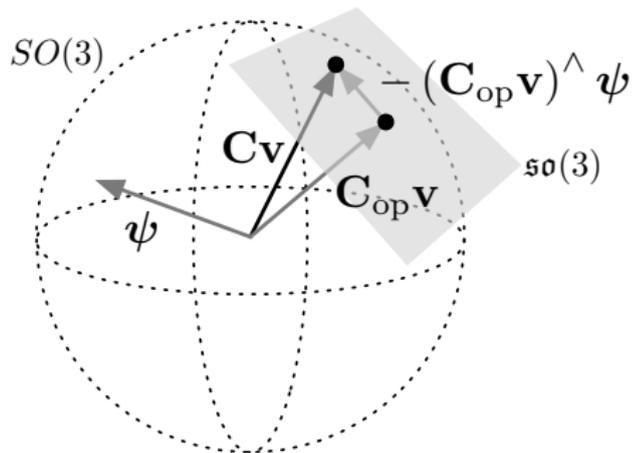
$$\mathbf{C}\mathbf{v} = \delta\mathbf{C}\bar{\mathbf{C}}\mathbf{v} = \exp(\psi^\wedge)\bar{\mathbf{C}}\mathbf{v} \approx (\mathbf{1} + \psi^\wedge)\bar{\mathbf{C}}\mathbf{v} = \bar{\mathbf{C}}\mathbf{v} - (\bar{\mathbf{C}}\mathbf{v})^\wedge\psi \quad (112)$$

Rotation perturbations



$$\mathbf{v} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \psi = \begin{bmatrix} 0 \\ 0 \\ \delta\theta \end{bmatrix} \quad \bar{\mathbf{C}} = \begin{bmatrix} \cos \bar{\theta} & -\sin \bar{\theta} & 0 \\ \sin \bar{\theta} & \cos \bar{\theta} & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Rotation perturbations



$$\mathbf{C}\mathbf{v} = \exp(\psi^\wedge) \mathbf{C}_{op}\mathbf{v} \approx \mathbf{C}_{op}\mathbf{v} - (\mathbf{C}_{op}\mathbf{v})^\wedge \psi$$

Pose perturbations

- for **poses**, we will perturb like this:

$$\mathbf{T} = \underbrace{\delta\mathbf{T}}_{\text{'small'}} \underbrace{\bar{\mathbf{T}}}_{\text{'big'}} \quad (113)$$

- we pick the following **perturbation**

$$\delta\mathbf{T} = \exp(\epsilon^\wedge) \quad (114)$$

which ensures \mathbf{T} is still a valid pose (by closure)

- can now **linearize** the product of a pose and homogeneous point, \mathbf{p} :

$$\mathbf{T}\mathbf{p} = \delta\mathbf{T}\bar{\mathbf{T}}\mathbf{p} = \exp(\epsilon^\wedge)\bar{\mathbf{T}}\mathbf{p} \approx (\mathbf{1} + \epsilon^\wedge)\bar{\mathbf{T}}\mathbf{p} = \bar{\mathbf{T}}\mathbf{p} + (\bar{\mathbf{T}}\mathbf{p})^\odot \epsilon \quad (115)$$

- we have used

$$\epsilon^\wedge \mathbf{p} \equiv \mathbf{p}^\odot \epsilon, \quad \mathbf{p}^\odot = \begin{bmatrix} \boldsymbol{\rho} \\ \eta \end{bmatrix}^\odot = \begin{bmatrix} \eta \mathbf{1} & -\boldsymbol{\rho}^\wedge \\ \mathbf{0}^T & \mathbf{0}^T \end{bmatrix} \quad (116)$$

Rotation optimization

- choose a perturbation scheme,

$$\mathbf{C} = \exp(\hat{\boldsymbol{\psi}}) \mathbf{C}_{\text{op}} \quad (117)$$

where $\boldsymbol{\psi}$ is a small **perturbation** applied to an **initial guess**, \mathbf{C}_{op}

- insert this in the function, $u(\mathbf{x})$, to be optimized:

$$\begin{aligned} u(\mathbf{C}\mathbf{v}) &= u(\exp(\hat{\boldsymbol{\psi}}) \mathbf{C}_{\text{op}}\mathbf{v}) \approx u((\mathbf{1} + \hat{\boldsymbol{\psi}}) \mathbf{C}_{\text{op}}\mathbf{v}) \\ &\approx u(\mathbf{C}_{\text{op}}\mathbf{v}) - \underbrace{\frac{\partial u}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{C}_{\text{op}}\mathbf{v}} (\mathbf{C}_{\text{op}}\mathbf{v})^{\wedge} \boldsymbol{\psi}}_{\boldsymbol{\delta}^T} = u(\mathbf{C}_{\text{op}}\mathbf{v}) + \boldsymbol{\delta}^T \boldsymbol{\psi} \end{aligned} \quad (118)$$

- then pick a perturbation, $\boldsymbol{\psi}$, to decrease the function

Rotation optimization: gradient descent

- suppose we would like to perform **Riemannian gradient descent**
- in this case, we would pick the perturbation to be of the form

$$\psi = -\alpha\delta \quad (119)$$

with $\alpha > 0$ a small step size

- we see the function is reduced by taking this step:

$$u(\mathbf{C}\mathbf{v}) - u(\mathbf{C}_{\text{op}}\mathbf{v}) \approx -\underbrace{\alpha\delta^T\delta}_{\geq 0} \quad (120)$$

- **retraction**: apply the perturbation to update the initial guess,

$$\mathbf{C}_{\text{op}} \leftarrow \exp(-\alpha\delta^\wedge) \mathbf{C}_{\text{op}} \quad (121)$$

so that $\mathbf{C}_{\text{op}} \in SO(3)$ at each iteration; iterate to convergence

Rotation optimization: Gauss-Newton

- gradient descent can be quite slow
- let's look at **Gauss-Newton** optimization
- suppose we have a general nonlinear, quadratic cost function of a rotation of the form,

$$J(\mathbf{C}) = \frac{1}{2} \sum_m (u_m(\mathbf{C}\mathbf{v}_m))^2 \quad (122)$$

where $u_m(\cdot)$ are scalar nonlinear functions and $\mathbf{v}_m \in \mathbb{R}^3$ are three-dimensional points

- we begin with an initial guess for the optimal rotation, $\mathbf{C}_{\text{op}} \in SO(3)$, and then perturb this (on the left) according to

$$\mathbf{C} = \exp(\boldsymbol{\psi}^\wedge) \mathbf{C}_{\text{op}} \quad (123)$$

where $\boldsymbol{\psi}$ is the perturbation

Rotation optimization: Gauss-Newton

- we then apply our perturbation scheme inside each $u_m(\cdot)$ so that

$$\begin{aligned} u_m(\mathbf{C}\mathbf{v}_m) &= u_m(\exp(\boldsymbol{\psi}^\wedge)\mathbf{C}_{\text{op}}\mathbf{v}_m) \approx u_m((\mathbf{1} + \boldsymbol{\psi}^\wedge)\mathbf{C}_{\text{op}}\mathbf{v}_m) \\ &\approx \underbrace{u_m(\mathbf{C}_{\text{op}}\mathbf{v}_m)}_{\beta_m} - \underbrace{\frac{\partial u_m}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{C}_{\text{op}}\mathbf{v}_m} (\mathbf{C}_{\text{op}}\mathbf{v}_m)^\wedge \boldsymbol{\psi}}_{\boldsymbol{\delta}_m^T} \end{aligned} \quad (124)$$

is a linearized version of $u_m(\cdot)$ in terms of our perturbation, $\boldsymbol{\psi}$

- inserting this back into our cost function we have

$$J(\mathbf{C}) \approx \frac{1}{2} \sum_m (\boldsymbol{\delta}_m^T \boldsymbol{\psi} + \beta_m)^2 \quad (125)$$

which is **exactly quadratic** in $\boldsymbol{\psi}$

Rotation optimization: Gauss-Newton

- taking the derivative of J with respect to $\boldsymbol{\psi}$ we have

$$\frac{\partial J}{\partial \boldsymbol{\psi}^T} = \sum_m \boldsymbol{\delta}_m (\boldsymbol{\delta}_m^T \boldsymbol{\psi} + \beta_m) \quad (126)$$

- set the derivative to zero to find the **optimal perturbation**, $\boldsymbol{\psi}^*$:

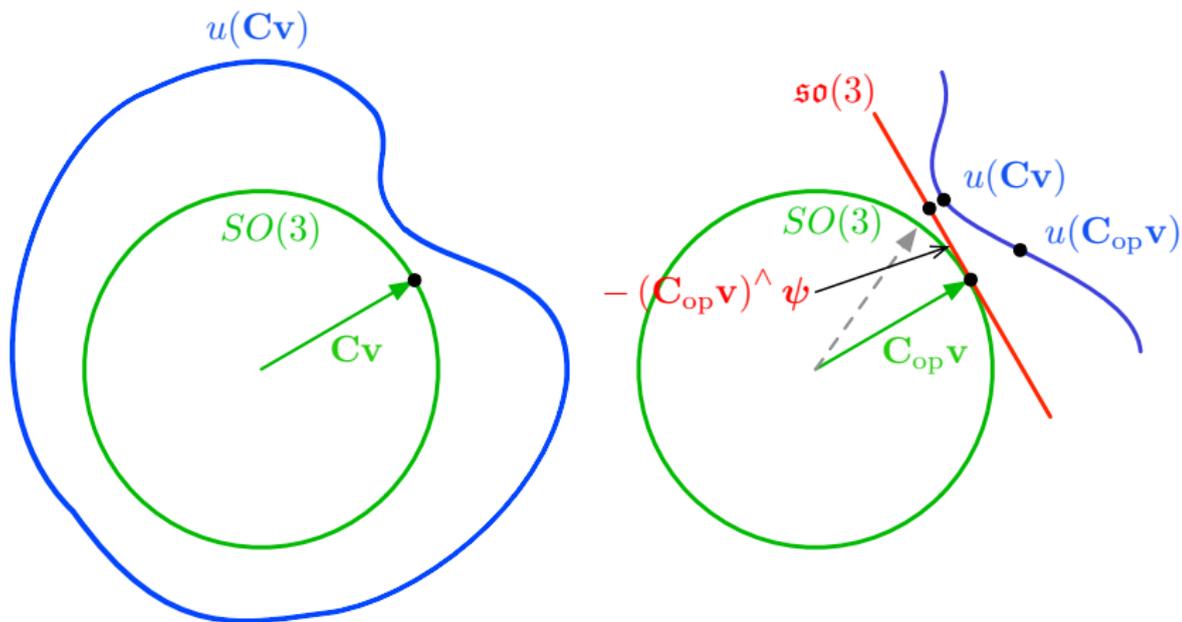
$$\left(\sum_m \boldsymbol{\delta}_m \boldsymbol{\delta}_m^T \right) \boldsymbol{\psi}^* = - \sum_m \beta_m \boldsymbol{\delta}_m \quad (127)$$

- this is a linear system of equations, which we can solve for $\boldsymbol{\psi}^*$
- **retraction**: apply this optimal perturbation to our initial guess,

$$\mathbf{C}_{\text{op}} \leftarrow \exp \left(\boldsymbol{\psi}^{*\wedge} \right) \mathbf{C}_{\text{op}}, \quad (128)$$

so that $\mathbf{C}_{\text{op}} \in SO(3)$ at each iteration; iterate to convergence

Rotation optimization



Rotation optimization commentary

- we have adapted classic optimization algorithms to work with the matrix Lie group, $SO(3)$, by exploiting the surjective property of the exponential map to define an appropriate perturbation scheme

$$\mathbf{C} = \exp(\psi^\wedge) \mathbf{C}_{\text{op}} \quad (129)$$

- we are essentially assuming that at each iteration the update, ψ , will be small and so have mapped the optimization problem from the Lie group up into the Lie algebra, $\mathfrak{so}(3)$
- this approach has three major advantages:
 - we are storing our rotation in a **singularity-free format**, \mathbf{C}_{op}
 - at each iteration we are performing **unconstrained optimization**
 - our manipulations occur at the **matrix level**
- we get away with this because the perturbation always becomes very small as we converge to the optimum

Pose optimization

- the same concepts can also be applied to poses
- suppose we have a general nonlinear, quadratic cost function of a transformation of the form

$$J(\mathbf{T}) = \frac{1}{2} \sum_m (u_m(\mathbf{T}\mathbf{p}_m))^2 \quad (130)$$

where $u_m(\cdot)$ are nonlinear functions and $\mathbf{p}_m \in \mathbb{R}^4$ are three-dimensional points expressed in **homogeneous coordinates**

- we begin with an initial guess for the optimal transformation, $\mathbf{T}_{\text{op}} \in SE(3)$, and then perturb this (on the left) according to

$$\mathbf{T} = \exp(\epsilon^\wedge) \mathbf{T}_{\text{op}} \quad (131)$$

where ϵ is the perturbation

Pose optimization

- we then apply our perturbation scheme inside each $u_m(\cdot)$ so that

$$\begin{aligned} u_m(\mathbf{T}\mathbf{p}_m) &= u_m(\exp(\boldsymbol{\epsilon}^\wedge)\mathbf{T}_{\text{op}}\mathbf{p}_m) \approx u_m((\mathbf{1} + \boldsymbol{\epsilon}^\wedge)\mathbf{T}_{\text{op}}\mathbf{p}_m) \\ &\approx \underbrace{u_m(\mathbf{T}_{\text{op}}\mathbf{p}_m)}_{\beta_m} + \underbrace{\frac{\partial u_m}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{T}_{\text{op}}\mathbf{p}_m} (\mathbf{T}_{\text{op}}\mathbf{p}_m)^\odot}_{\boldsymbol{\delta}_m^T} \boldsymbol{\epsilon} \end{aligned} \quad (132)$$

is a linearized version of $u_m(\cdot)$ in terms of our perturbation, $\boldsymbol{\epsilon}$

- inserting this back into our cost function we have

$$J(\mathbf{T}) = \frac{1}{2} \sum_m (\boldsymbol{\delta}_m^T \boldsymbol{\epsilon} + \beta_m)^2 \quad (133)$$

which is exactly quadratic in $\boldsymbol{\epsilon}$

Pose optimization

- taking the derivative of J with respect to ϵ we have

$$\frac{\partial J}{\partial \epsilon^T} = \sum_m \delta_m (\delta_m^T \epsilon + \beta_m) \quad (134)$$

- set the derivative to zero to find the **optimal perturbation, ϵ^*** :

$$\left(\sum_m \delta_m \delta_m^T \right) \epsilon^* = - \sum_m \beta_m \delta_m \quad (135)$$

- this is a linear system of equations, which we can solve for ϵ^*
- **retraction**: apply this optimal perturbation to our initial guess,

$$\mathbf{T}_{\text{op}} \leftarrow \exp(\epsilon^{*\wedge}) \mathbf{T}_{\text{op}} \quad (136)$$

so that $\mathbf{T}_{\text{op}} \in SE(3)$ at each iteration; iterate to convergence

Point-cloud alignment

- consider the problem of **aligning two point-clouds**, \mathbf{y}_j and \mathbf{p}_j , which are in homogeneous-point form and $j = 1 \dots J$
- we define our error term for each point pair as

$$\mathbf{e}_j = \mathbf{y}_j - \mathbf{T}\mathbf{p}_j \quad (137)$$

- we define our objective function as

$$J(\mathbf{T}) = \frac{1}{2} \sum_{j=1}^M w_j \mathbf{e}_j^T \mathbf{e}_j = \frac{1}{2} \sum_{j=1}^M w_j (\mathbf{y}_j - \mathbf{T}\mathbf{p}_j)^T (\mathbf{y}_j - \mathbf{T}\mathbf{p}_j) \quad (138)$$

where $w_j > 0$ are scalar weights

- we seek to minimize J with respect to $\mathbf{T} \in SE(3)$; we want to know the pose between the two point-clouds

Point-cloud alignment

- we use our $SE(3)$ -sensitive perturbation scheme

$$\mathbf{T} = \exp(\epsilon^\wedge) \mathbf{T}_{\text{op}} \approx (\mathbf{1} + \epsilon^\wedge) \mathbf{T}_{\text{op}} \quad (139)$$

where \mathbf{T}_{op} is some initial guess and ϵ is a small perturbation

- inserting this into the objective function we then have

$$J(\mathbf{T}) \approx \frac{1}{2} \sum_{j=1}^M w_j \left((\mathbf{y}_j - \mathbf{z}_j) - \mathbf{z}_j^\odot \epsilon \right)^T \left((\mathbf{y}_j - \mathbf{z}_j) - \mathbf{z}_j^\odot \epsilon \right) \quad (140)$$

where $\mathbf{z}_j = \mathbf{T}_{\text{op}} \mathbf{p}_j$ and we have used that

$$\epsilon^\wedge \mathbf{z}_j = \mathbf{z}_j^\odot \epsilon \quad (141)$$

- the objective function is now **exactly quadratic** in ϵ

Point-cloud alignment

- we can carry out a simple, **unconstrained optimization** for ϵ
- taking the derivative we find

$$\frac{\partial J}{\partial \epsilon^T} = - \sum_{j=1}^M w_j \mathbf{z}_j^{\odot T} \left((\mathbf{y}_j - \mathbf{z}_j) - \mathbf{z}_j^{\odot} \epsilon \right) \quad (142)$$

- setting this to zero, we have the following system of equations for the **optimal ϵ^*** :

$$\left(\frac{1}{w} \sum_{j=1}^M w_j \mathbf{z}_j^{\odot T} \mathbf{z}_j^{\odot} \right) \epsilon^* = \frac{1}{w} \sum_{j=1}^M w_j \mathbf{z}_j^{\odot T} (\mathbf{y}_j - \mathbf{z}_j) \quad (143)$$

- **retraction**: apply the optimal perturbation and iterate to convergence:

$$\mathbf{T}_{\text{op}} \leftarrow \exp \left(\epsilon^{\wedge} \right) \mathbf{T}_{\text{op}} \quad (144)$$

Pose optimization commentary

- we have adapted classic optimization algorithms to work with the matrix Lie group, $SE(3)$, by exploiting the surjective property of the exponential map to define an appropriate perturbation scheme

$$\mathbf{T} = \exp(\epsilon^\wedge) \mathbf{T}_{\text{op}} \quad (145)$$

- we are essentially assuming that at each iteration the update, ϵ , will be small and so have mapped the optimization problem from the Lie group up into the Lie algebra, $\mathfrak{se}(3)$
- this approach has three major advantages:
 - we are storing our pose in a **singularity-free format**, \mathbf{T}_{op}
 - at each iteration we are performing **unconstrained optimization**
 - our manipulations occur at the **matrix level**
- we get away with this because the perturbation always becomes very small as we converge to the optimum

Lecture Summary

- we have seen that despite **nonlinearities** in our motion and observation models and even state variables that are members of Lie groups, we can still carry out **MAP** estimation very similar to the linear-Gaussian case
- the difference is that we now need to **iterate** to convergence, which we do using some variant of **Gauss-Newton**
- many of the usual algorithms differ only based on which variables are held fixed and which variables are optimized at each iteration (e.g., full batch, extended RTS smoother, sliding-window filter, (iterated) EKF)
- **Lecture 3**: Continuous-Time Estimation

Lecture 3: Continuous-Time Estimation

A Short Course in State Estimation

Timothy D. Barfoot

University of Toronto

Copyright © 2024

Lecture Outline

Lecture 1: Linear-Gaussian Estimation

Lecture 2: Nonlinearities and Lie Groups

Lecture 3: Continuous-Time Estimation

Problem Setup

Motion Prior

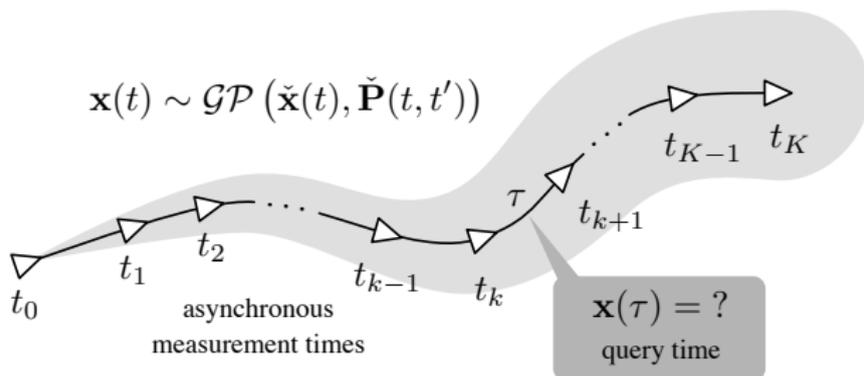
GP Regression

Lie Groups

STEAM

Lecture 4: Beyond MAP – Variational Inference

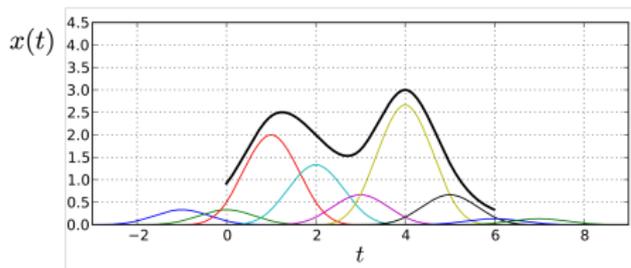
Motivation



- there are situations where it is beneficial to think of the robot's trajectory as a **continuous function of time**
- e.g., asynchronous measurements, scanning-while-moving sensors

Two camps

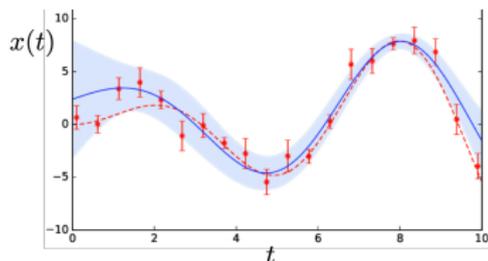
$$\mathbf{x}(t) = \Psi(t)\mathbf{c}$$



Parametric methods:

- (i) hard to decide what basis functions and how many
- (ii) don't encode motion priors explicitly
- (iii) computations are typically fast

$$\mathbf{x}(t) \sim \mathcal{GP}(\check{\mathbf{x}}(t), \check{\mathbf{P}}(t, t'))$$



Nonparametric methods:

- (i) no need to pick number of basis functions
- (ii) can encode motion priors explicitly
- (iii) can be computationally expensive

- two main camps: **parametric** (e.g., spline) and **nonparametric** (e.g., Gaussian process)
- we will explore the latter in this lecture

System

- we define our system using the following **linear, time-varying** models:

motion model: $\dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{v}(t) + \mathbf{L}(t)\mathbf{w}(t)$ (146a)

observation model: $\mathbf{y}(t_k) = \mathbf{C}(t_k)\mathbf{x}(t_k) + \mathbf{n}(t_k), \quad k = 0 \dots K$ (146b)

where k is again the discrete-time index and K its maximum

- the variables have the following meanings:

system state : $\mathbf{x}(t) \in \mathbb{R}^N$

initial state : $\mathbf{x}_0 \in \mathbb{R}^N \sim \mathcal{N}(\check{\mathbf{x}}_0, \check{\mathbf{P}}_0)$

input : $\mathbf{v}(t) \in \mathbb{R}^N$

process noise : $\mathbf{w}(t) \in \mathbb{R}^N \sim \mathcal{GP}(\mathbf{0}, \mathbf{Q}\delta(t-t'))$

measurement : $\mathbf{y}(t_k) \in \mathbb{R}^M$

measurement noise : $\mathbf{n}(t_k) \in \mathbb{R}^M \sim \mathcal{N}(\mathbf{0}, \mathbf{R}(t_k))$

Motion prior

- to connect to **Gaussian process regression**, we will convert our motion model

$$\dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{v}(t) + \mathbf{L}(t)\mathbf{w}(t) \quad (147)$$

into a **kernel function**

- after **stochastic integration** we have a **GP motion prior**

$$\mathbf{x}(t) \sim \mathcal{GP} \left(\underbrace{\Phi(t, t_0)\check{\mathbf{x}}_0 + \int_{t_0}^t \Phi(t, s)\mathbf{v}(s) ds}_{\check{\mathbf{x}}(t)}, \underbrace{\Phi(t, t_0)\check{\mathbf{P}}_0\Phi(t', t_0)^T + \int_{t_0}^{\min(t, t')} \Phi(t, s)\mathbf{L}(s)\mathbf{Q}\mathbf{L}(s)^T\Phi(t', s)^T ds}_{\check{\mathbf{P}}(t, t')} \right) \quad (148)$$

where $\Phi(t, s)$ is known as the **transition function**

Estimate at measurement times

- if we only want to solve at the **measurement times**, things are basically identical to the **linear-Gaussian** case:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} (\underbrace{\mathbf{A}\mathbf{v}}_{\check{\mathbf{p}}} - \mathbf{x})^T \underbrace{\mathbf{A}^{-T}\mathbf{Q}^{-1}\mathbf{A}^{-1}}_{\check{\mathbf{P}}^{-1}} (\underbrace{\mathbf{A}\mathbf{v}}_{\check{\mathbf{p}}} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{C}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{C}\mathbf{x}) \quad (150)$$

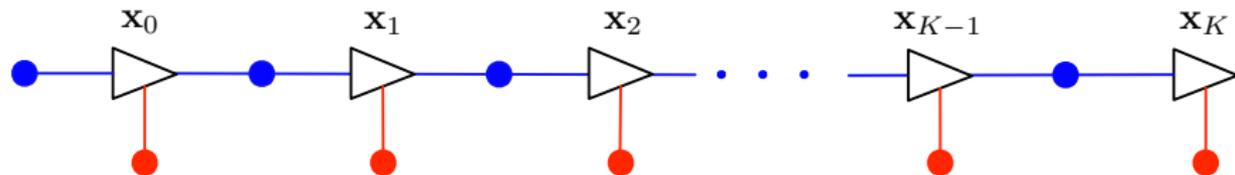
$$\Rightarrow \underbrace{(\mathbf{A}^{-T}\mathbf{Q}^{-1}\mathbf{A}^{-1} + \mathbf{C}^T\mathbf{R}^{-1}\mathbf{C})}_{\text{block-tridiagonal}} \hat{\mathbf{x}} = \mathbf{A}^{-T}\mathbf{Q}^{-1}\mathbf{v} + \mathbf{C}^T\mathbf{R}^{-1}\mathbf{y} \quad (151)$$

where

$$\hat{\mathbf{x}} = \begin{bmatrix} \hat{\mathbf{x}}(t_0) \\ \vdots \\ \hat{\mathbf{x}}(t_K) \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \mathbf{y}(t_0) \\ \vdots \\ \mathbf{y}(t_K) \end{bmatrix}, \quad \mathbf{C} = \text{diag}(\mathbf{C}(t_0), \dots, \mathbf{C}(t_K)), \quad \mathbf{R} = \text{diag}(\mathbf{R}(t_0), \dots, \mathbf{R}(t_K))$$

- only now our motion prior has been constructed from a **stochastic differential equation**

Still a sparse factor graph



$$\underbrace{(\underbrace{A^{-T}Q^{-1}A^{-1}}_{\text{motion prior}} + \underbrace{C^TR^{-1}C}_{\text{measurements}})}_{\hat{\mathbf{P}}^{-1}} \hat{\mathbf{x}} = A^{-T}Q^{-1}\mathbf{v} + C^TR^{-1}\mathbf{y}$$

block-tridiagonal

Querying at other times

- suppose that we'd like to **query** the trajectory as some other times $\tau_0 < \tau_1 < \dots < \tau_J$
- the **joint density** between the state (at the query times) and the measurements (at the measurement times) is written as

$$p\left(\begin{bmatrix} \mathbf{x}_\tau \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \check{\mathbf{x}}_\tau \\ \mathbf{C}\check{\mathbf{x}} \end{bmatrix}, \begin{bmatrix} \check{\mathbf{P}}_{\tau\tau} & \check{\mathbf{P}}_\tau \mathbf{C}^T \\ \mathbf{C}\check{\mathbf{P}}_\tau^T & \mathbf{R} + \mathbf{C}\check{\mathbf{P}}\mathbf{C}^T \end{bmatrix}\right) \quad (152)$$

where

$$\check{\mathbf{x}} = \begin{bmatrix} \check{\mathbf{x}}(t_0) \\ \vdots \\ \check{\mathbf{x}}(t_K) \end{bmatrix}, \quad \mathbf{x}_\tau = \begin{bmatrix} \mathbf{x}(\tau_0) \\ \vdots \\ \mathbf{x}(\tau_J) \end{bmatrix}, \quad \check{\mathbf{x}}_\tau = \begin{bmatrix} \check{\mathbf{x}}(\tau_0) \\ \vdots \\ \check{\mathbf{x}}(\tau_J) \end{bmatrix}, \quad \check{\mathbf{P}}_\tau = [\check{\mathbf{P}}(\tau_i, t_j)]_{ij}, \quad \check{\mathbf{P}}_{\tau\tau} = [\check{\mathbf{P}}(\tau_i, \tau_j)]_{ij}$$

- using the **Gaussian conditioning** formula we have

$$p(\mathbf{x}_\tau | \mathbf{y}) = \mathcal{N}\left(\underbrace{\check{\mathbf{x}}_\tau + \check{\mathbf{P}}_\tau \mathbf{C}^T (\mathbf{C}\check{\mathbf{P}}\mathbf{C}^T + \mathbf{R})^{-1} (\mathbf{y} - \mathbf{C}\check{\mathbf{x}})}_{\check{\mathbf{x}}_\tau, \text{ mean}}, \underbrace{\check{\mathbf{P}}_{\tau\tau} - \check{\mathbf{P}}_\tau \mathbf{C}^T (\mathbf{C}\check{\mathbf{P}}\mathbf{C}^T + \mathbf{R})^{-1} \mathbf{C}\check{\mathbf{P}}_\tau^T}_{\check{\mathbf{P}}_{\tau\tau}, \text{ covariance}}\right) \quad (153)$$

Querying at other times

- since we have already solved at the measurement times, we can use $\hat{\mathbf{x}}$ and $\hat{\mathbf{P}}$ to simplify the solution at the query times:

$$\hat{\mathbf{x}}_\tau = \check{\mathbf{x}}_\tau + (\check{\mathbf{P}}_\tau \check{\mathbf{P}}^{-1}) (\hat{\mathbf{x}} - \check{\mathbf{x}}) \quad (154a)$$

$$\hat{\mathbf{P}}_{\tau\tau} = \check{\mathbf{P}}_{\tau\tau} + (\check{\mathbf{P}}_\tau \check{\mathbf{P}}^{-1}) (\hat{\mathbf{P}} - \check{\mathbf{P}}) (\check{\mathbf{P}}_\tau \check{\mathbf{P}}^{-1})^T \quad (154b)$$

which are the **standard GP interpolation equations**

- our inverse kernel matrix, $\check{\mathbf{P}}^{-1}$, is **block-tridiagonal**, so interpolation is **extremely efficient** compared to generic GP regression

Single query is $O(1)$

- consider a **single query time**, $t_k \leq \tau < t_{k+1}$
- we can show that

$$\hat{\mathbf{x}}(\tau) = \check{\mathbf{x}}(\tau) + [\mathbf{\Lambda}(\tau) \quad \mathbf{\Psi}(\tau)] \left(\begin{bmatrix} \hat{\mathbf{x}}_k \\ \hat{\mathbf{x}}_{k+1} \end{bmatrix} - \begin{bmatrix} \check{\mathbf{x}}(t_k) \\ \check{\mathbf{x}}(t_{k+1}) \end{bmatrix} \right) \quad (155a)$$

$$\hat{\mathbf{P}}(\tau, \tau) = \check{\mathbf{P}}(\tau, \tau) + [\mathbf{\Lambda}(\tau) \quad \mathbf{\Psi}(\tau)] \left(\begin{bmatrix} \hat{\mathbf{P}}_{k,k} & \hat{\mathbf{P}}_{k,k+1} \\ \hat{\mathbf{P}}_{k+1,k} & \hat{\mathbf{P}}_{k+1,k+1} \end{bmatrix} - \begin{bmatrix} \check{\mathbf{P}}(t_k, t_k) & \check{\mathbf{P}}(t_k, t_{k+1}) \\ \check{\mathbf{P}}(t_{k+1}, t_k) & \check{\mathbf{P}}(t_{k+1}, t_{k+1}) \end{bmatrix} \right) \begin{bmatrix} \mathbf{\Lambda}(\tau)^T \\ \mathbf{\Psi}(\tau)^T \end{bmatrix} \quad (155b)$$

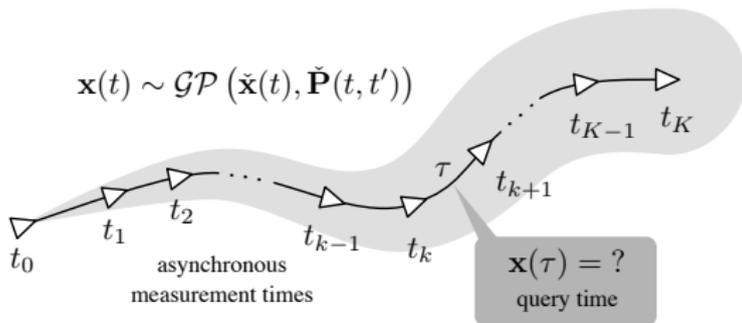
where

$$\mathbf{\Lambda}(\tau) = \mathbf{\Phi}(\tau, t_k) - \mathbf{Q}_\tau \mathbf{\Phi}(t_{k+1}, \tau)^T \mathbf{Q}_{k+1}^{-1} \mathbf{\Phi}(t_{k+1}, t_k), \quad \mathbf{\Psi}(\tau) = \mathbf{Q}_\tau \mathbf{\Phi}(t_{k+1}, \tau)^T \mathbf{Q}_{k+1}^{-1}, \quad (156)$$

$$\mathbf{Q}_\tau = \int_{t_k}^{\tau} \mathbf{\Phi}(\tau, s) \mathbf{L}(s) \mathbf{Q} \mathbf{L}(s)^T \mathbf{\Phi}(\tau, s)^T ds \quad (157)$$

- in other words, we only need to use the posterior solution at the **two measurement times that bracket** the query time

GP regression summary



- solve at the **measurement times** $t_0 < t_1 < \dots < t_K$ in $O(K)$ time
- solve at each **query time**, τ , in $O(1)$ time
- principled approach to interpolation since we're using the motion model not some other scheme
- we can also use this GP interpolation to reduce the number of state variables while accounting for all measurement timestamps – **GP inducing points**

Example: white-noise on acceleration prior

- consider the case $\ddot{\mathbf{p}}(t) = \mathbf{w}(t)$, where $\mathbf{p}(t)$ corresponds to position and $\mathbf{w}(t) \sim \mathcal{GP}(\mathbf{0}, \mathbf{Q} \delta(t - t'))$ is white noise as before
- we can cast this in the form

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{v}(t) + \mathbf{L}\mathbf{w}(t) \quad (158)$$

by taking

$$\mathbf{x}(t) = \begin{bmatrix} \mathbf{p}(t) \\ \dot{\mathbf{p}}(t) \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{v}(t) = \mathbf{0}, \quad \mathbf{L} = \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix} \quad (159)$$

- importantly, the **Markovian state** is now both position and velocity
- in the absence of measurements to the contrary, this motion prior tries to keep the velocity constant

Example: white-noise on acceleration prior

- the transition function for this **linear, time-invariant** example is

$$\Phi(t_k, t_{k-1}) = \begin{bmatrix} \mathbf{1} & \Delta t_{k:k-1} \mathbf{1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (160)$$

where $\Delta t_{k:k-1} = t_k - t_{k-1}$

- since there is no input, $\mathbf{v}(t) = \mathbf{0}$, the prior becomes

$$\check{\mathbf{x}} = \mathbf{A}\mathbf{v}, \quad \check{\mathbf{P}} = \mathbf{A}\mathbf{Q}\mathbf{A}^T, \quad (161)$$

$$\mathbf{v} = \begin{bmatrix} \check{\mathbf{x}}_0 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{Q} = \text{diag}(\check{\mathbf{P}}_0, \mathbf{Q}_1, \dots, \mathbf{Q}_K), \quad \mathbf{Q}_k = \begin{bmatrix} \frac{1}{3}\Delta t_{k:k-1}^3 \mathbf{Q} & \frac{1}{2}\Delta t_{k:k-1}^2 \mathbf{Q} \\ \frac{1}{2}\Delta t_{k:k-1}^2 \mathbf{Q} & \Delta t_{k:k-1} \mathbf{Q} \end{bmatrix} \quad (162)$$

Example: white-noise on acceleration prior

- after solving at the measurement times, the interpolation equations for $t_k \leq \tau < t_{k+1}$ become

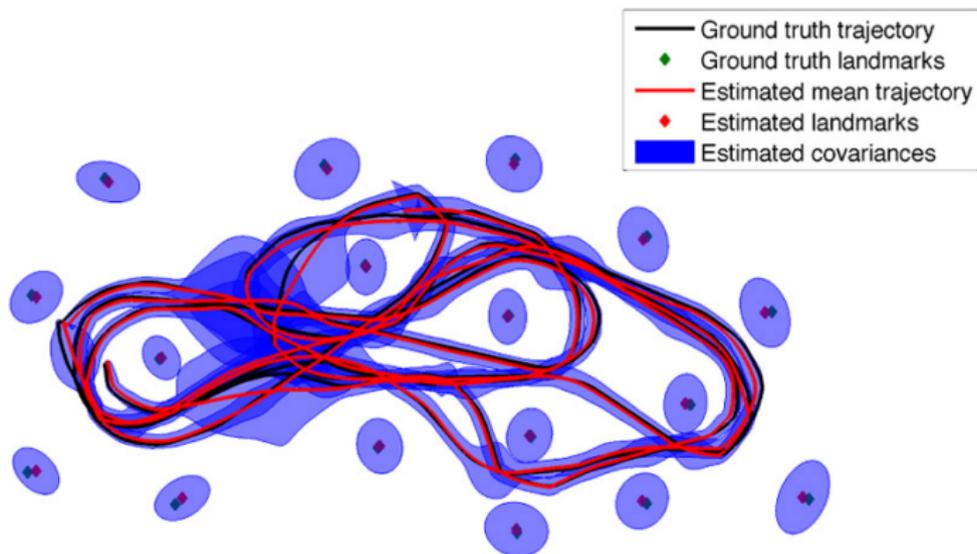
$$\hat{\mathbf{x}}_\tau = \check{\mathbf{x}}_\tau + \begin{bmatrix} (1 - 3\alpha^2 + 2\alpha^3)\mathbf{1} & T(\alpha - 2\alpha^2 + \alpha^3)\mathbf{1} \\ \frac{1}{T}6(-\alpha + \alpha^2)\mathbf{1} & (1 - 4\alpha + 3\alpha^2)\mathbf{1} \end{bmatrix} (\hat{\mathbf{x}}_k - \check{\mathbf{x}}_k) + \begin{bmatrix} (3\alpha^2 - 2\alpha^3)\mathbf{1} & T(-\alpha^2 + \alpha^3)\mathbf{1} \\ \frac{1}{T}6(\alpha - \alpha^2)\mathbf{1} & (-2\alpha + 3\alpha^2)\mathbf{1} \end{bmatrix} (\hat{\mathbf{x}}_{k+1} - \check{\mathbf{x}}_{k+1}) \quad (163)$$

where

$$\check{\mathbf{x}}_\tau = \Phi(\tau, t_k)\check{\mathbf{x}}_k, \quad \alpha = \frac{\tau - t_k}{t_{k+1} - t_k} \in [0, 1], \quad T = \Delta t_{k+1:k} = t_{k+1} - t_k$$

- remarkably, the top row (corresponding to position) is precisely a **cubic Hermite polynomial** interpolation – bottom row (corresponding to velocity) is derivative of this
- this falls out of the math for free simply by picking the **WNOA motion prior**

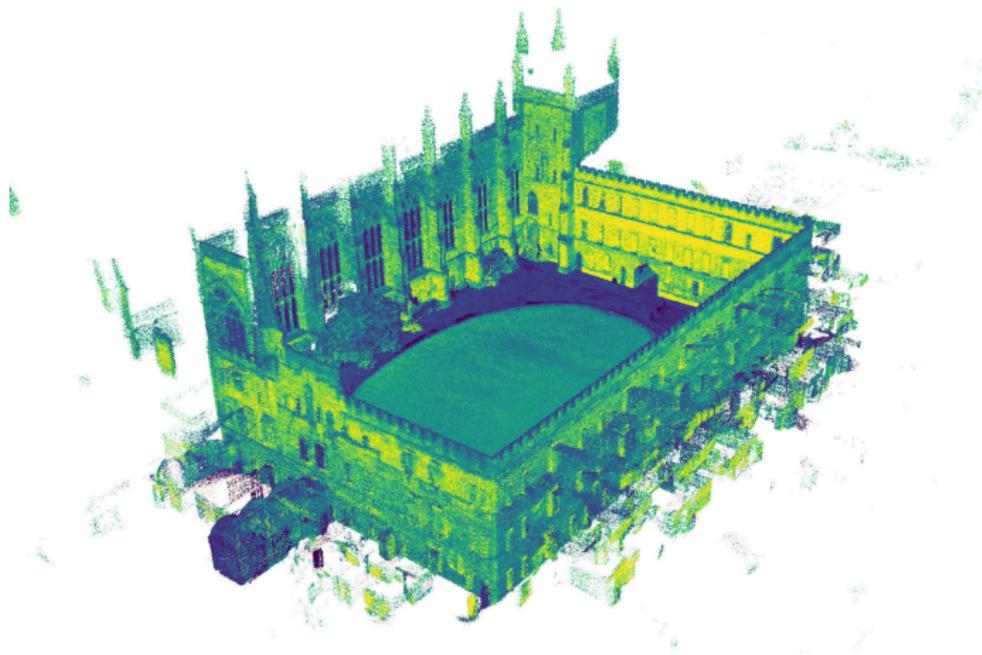
Example: white-noise on acceleration prior



- example **continuous-time** estimation using only the **WNOA motion prior** plus landmark measurements – no odometry!

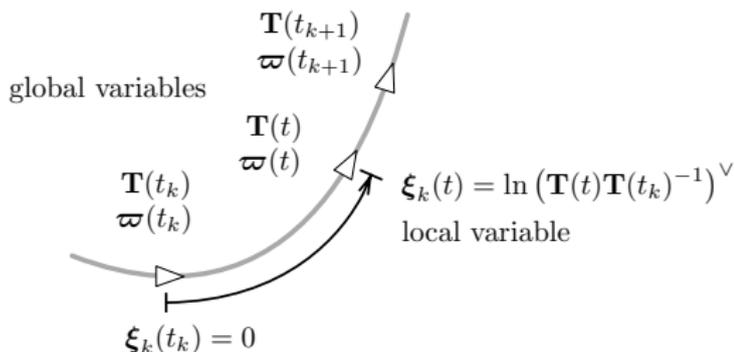
What about 3D?

- we would like to do **continuous-time** estimation with **Lie-group** state variables so we can estimate 3D trajectories



Special Euclidean group

- in order to recycle our **linear continuous-time** tools, we will stitch together many **local GPs**



- we introduce **local variables**, $\xi_k(t)$, in the Lie algebra of $SE(3)$ and then define a **WNOA motion prior** over these

$$\ddot{\xi}_k(t) = \mathbf{w}_k(t), \quad \mathbf{w}_k(t) \sim \mathcal{GP}(\mathbf{0}, \mathbf{Q}(t - t')) \quad (164)$$

Motion prior in Lie algebra

- we reorganize the **WNOA motion prior** into the requisite form

$$\frac{d}{dt} \begin{bmatrix} \boldsymbol{\xi}_k(t) \\ \boldsymbol{\psi}_k(t) \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \underbrace{\begin{bmatrix} \boldsymbol{\xi}_k(t) \\ \boldsymbol{\psi}_k(t) \end{bmatrix}}_{\boldsymbol{\gamma}_k(t)} + \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix} \mathbf{w}_k(t) \quad (165)$$

- we integrate this stochastically to obtain our **kernel function**

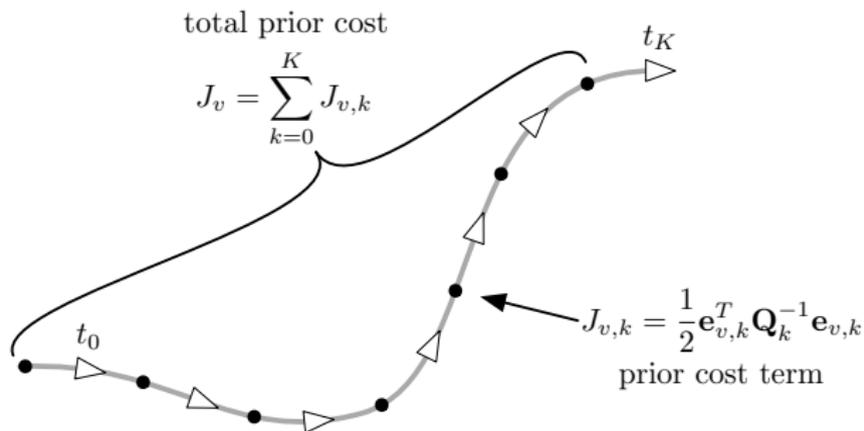
$$\boldsymbol{\gamma}_k(t) \sim \mathcal{GP} \left(\underbrace{\boldsymbol{\Phi}(t, t_k) \check{\boldsymbol{\gamma}}_k(t_k)}_{\text{mean function}}, \underbrace{\boldsymbol{\Phi}(t, t_k) \check{\mathbf{P}}(t_k) \boldsymbol{\Phi}(t, t_k)^T + \mathbf{Q}(t - t_k)}_{\text{covariance function}} \right) \quad (166)$$

where

$$\boldsymbol{\Phi}(t, t') = \begin{bmatrix} \mathbf{1} & (t - t')\mathbf{1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}, \quad \mathbf{Q}(t - t') = \begin{bmatrix} \frac{1}{3}(t - t')^3 \mathbf{Q} & \frac{1}{2}(t - t')^2 \mathbf{Q} \\ \frac{1}{2}(t - t')^2 \mathbf{Q} & (t - t') \mathbf{Q} \end{bmatrix}, \quad t \geq t' \quad (167)$$

and $\check{\boldsymbol{\gamma}}_k(t_k)$ and $\check{\mathbf{P}}(t_k)$ are the initial mean and covariance at $t = t_k$

Motion prior in Lie algebra

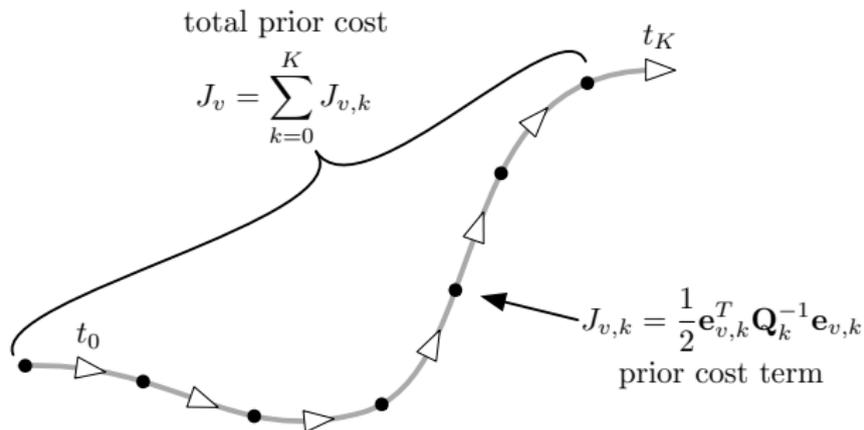


- build **errors** and assemble them into **MAP cost terms**

$$\mathbf{e}_{v,k} = \begin{cases} \begin{bmatrix} -\ln(\mathbf{T}(t_0)\check{\mathbf{T}}_0^{-1})^\vee \\ \check{\boldsymbol{\omega}}_0 - \boldsymbol{\omega}(t_0) \end{bmatrix} & k = 0 \\ \Phi(t_k, t_{k-1}) \begin{bmatrix} \gamma_{k-1}(t_{k-1}) - \check{\gamma}_{k-1}(t_{k-1}) \\ -(\gamma_{k-1}(t_k) - \check{\gamma}_{k-1}(t_k)) \end{bmatrix} & k > 0 \end{cases} \quad (168)$$

- note $\mathbf{Q}_0 = \check{\mathbf{P}}_0$ and $\mathbf{Q}_k = \mathbf{Q}(t_k - t_{k-1})$ for $k = 1 \dots K$

Motion prior in Lie algebra



- we back-substitute for the errors in terms of the **global variables**

$$\mathbf{e}_{v,k} = \begin{cases} \begin{bmatrix} -\ln(\mathbf{T}(t_0)\check{\mathbf{T}}_0^{-1})^\vee \\ \check{\boldsymbol{\omega}}_0 - \boldsymbol{\omega}(t_0) \end{bmatrix} & k = 0 \\ \begin{bmatrix} (t_k - t_{k-1}) \boldsymbol{\omega}(t_{k-1}) - \ln(\mathbf{T}(t_k)\mathbf{T}(t_{k-1})^{-1})^\vee \\ \boldsymbol{\omega}(t_{k-1}) - \mathcal{J}(\ln(\mathbf{T}(t_k)\mathbf{T}(t_{k-1})^{-1})^\vee)^{-1} \boldsymbol{\omega}(t_k) \end{bmatrix} & k > 0 \end{cases} \quad (169)$$

Lie group aside

- we need a few more **Lie-group tools** at this point
- the **$SE(3)$ Jacobian**, $\mathcal{J} \in \mathbb{R}^{6 \times 6}$, allows us to make the following approximation when ξ_1 small

$$\exp(\xi_1^\wedge) \exp(\xi_2^\wedge) \approx \exp(\mathcal{J}(\xi_2)^{-1} \xi_1 + \xi_2) \quad (170)$$

- the **adjoint** of a pose, $\mathcal{T} = \text{Ad}(\mathbf{T}) \in \mathbb{R}^{6 \times 6}$, is an **inner automorphism** operator that allows

$$(\mathcal{T}\mathbf{x})^\wedge = \mathbf{T}\mathbf{x}^\wedge\mathbf{T}^{-1} \quad (171)$$

Motion prior in Lie algebra

- with these **Lie group tools** we can perturb our to-be-estimated state variables according to

$$\mathbf{T}_k = \exp(\hat{\boldsymbol{\epsilon}}_k) \mathbf{T}_{\text{op},k}, \quad \boldsymbol{\varpi}_k = \boldsymbol{\varpi}_{\text{op},k} + \boldsymbol{\eta}_k \quad (172)$$

- we then **linearize** our motion prior errors terms so that

$$\mathbf{e}_{v,k}(\mathbf{x}) \approx \mathbf{e}_{v,k}(\mathbf{x}_{\text{op}}) + \mathbf{F}_{k-1} \boldsymbol{\varepsilon}_{k-1} - \mathbf{E}_k \boldsymbol{\varepsilon}_k \quad (173)$$

where $\boldsymbol{\varepsilon}_k = \begin{bmatrix} \boldsymbol{\epsilon}_k \\ \boldsymbol{\eta}_k \end{bmatrix}$ and

$$\mathbf{F}_{k-1} = \begin{cases} \begin{bmatrix} \mathbf{0} & (t_k - t_{k-1}) \mathbf{1} \\ \frac{1}{2} \boldsymbol{\varpi}_{\text{op},k}^\wedge \mathcal{J}_{\text{op},k,k-1}^{-1} \mathcal{T}_{\text{op},k,k-1} & \mathbf{1} \end{bmatrix} & k = 0 \\ \begin{bmatrix} \mathbf{0} & (t_k - t_{k-1}) \mathbf{1} \\ \frac{1}{2} \boldsymbol{\varpi}_{\text{op},k}^\wedge \mathcal{J}_{\text{op},k,k-1}^{-1} \mathcal{T}_{\text{op},k,k-1} & \mathbf{1} \end{bmatrix} & k > 0 \end{cases}, \quad \mathbf{E}_k = \begin{cases} \begin{bmatrix} \mathcal{J}(\ln(\mathbf{T}_{\text{op},0} \tilde{\mathbf{T}}_0^{-1}))^\vee^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} & k = 0 \\ \begin{bmatrix} \mathcal{J}_{\text{op},k,k-1}^{-1} & \mathbf{0} \\ \frac{1}{2} \boldsymbol{\varpi}_{\text{op},k}^\wedge \mathcal{J}_{\text{op},k,k-1}^{-1} & \mathcal{J}_{\text{op},k,k-1}^{-1} \end{bmatrix} & k > 0 \end{cases}$$

Motion prior in Lie algebra

- our **quadratic approximation** of the motion prior term at our current guess is therefore

$$J_v \approx \frac{1}{2} (\mathbf{e}_v(\mathbf{x}_{\text{op}}) - \mathbf{F}^{-1} \delta \mathbf{x}_1)^T \mathbf{Q}^{-1} (\mathbf{e}_v(\mathbf{x}_{\text{op}}) - \mathbf{F}^{-1} \delta \mathbf{x}_1) \quad (174)$$

where

$$\mathbf{e}_v(\mathbf{x}_{\text{op}}) = \begin{bmatrix} \mathbf{e}_{v,0}(\mathbf{x}_{\text{op}}) \\ \mathbf{e}_{v,1}(\mathbf{x}_{\text{op}}) \\ \vdots \\ \mathbf{e}_{v,K}(\mathbf{x}_{\text{op}}) \end{bmatrix}, \quad \mathbf{F}^{-1} = \begin{bmatrix} \mathbf{E}_0 & & & & & \\ -\mathbf{F}_0 & \mathbf{E}_1 & & & & \\ & -\mathbf{F}_1 & \mathbf{E}_2 & & & \\ & & \ddots & \ddots & & \\ & & & -\mathbf{F}_{K-1} & \mathbf{E}_K & \end{bmatrix}, \quad \delta \mathbf{x}_1 = \begin{bmatrix} \varepsilon_0 \\ \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_K \end{bmatrix},$$
$$\mathbf{Q} = \text{diag}(\tilde{\mathbf{P}}_0, \mathbf{Q}_1, \dots, \mathbf{Q}_K)$$

- this is essentially all we need to carry out **MAP estimation** via Gauss-Newton once we incorporate some measurement terms

STEAM

- let's look at **simultaneous trajectory estimation and mapping**
- the state to be estimated is now

$$\mathbf{x} = \left\{ \underbrace{\mathbf{T}_0, \boldsymbol{\varpi}_0, \dots, \mathbf{T}_K, \boldsymbol{\varpi}_K}_{\mathbf{x}_1}, \underbrace{\mathbf{p}_1, \dots, \mathbf{p}_M}_{\mathbf{x}_2} \right\} \quad (175)$$

which includes poses, generalized velocities, and landmarks, \mathbf{p}_j

- the measurements at times t_0, t_1, \dots, t_K are

$$\mathbf{y} = \{\mathbf{y}_{10}, \dots, \mathbf{y}_{M0}, \dots, \mathbf{y}_{1K}, \dots, \mathbf{y}_{MK}\} \quad (176)$$

- each measurement is of the form

$$\mathbf{y}_{jk} = \mathbf{s}(\mathbf{T}_k \mathbf{p}_j) + \mathbf{n}_{jk}, \quad \mathbf{n}_{jk} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{jk}) \quad (177)$$

where $\mathbf{s}(\cdot)$ is a nonlinear **sensor model** (e.g., a stereo camera)

Landmarks and measurement error term

- our landmarks are in **homogeneous coordinates** so we perturb according to

$$\mathbf{p}_j = \mathbf{p}_{\text{op},j} + \mathbf{D}\boldsymbol{\zeta}_j, \quad \mathbf{D} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad (178)$$

where $\mathbf{p}_{\text{op},j}$ is the operating point and $\boldsymbol{\zeta}_j$ is the perturbation

- the error associated with a measurement is then

$$\mathbf{e}_{y,jk}(\mathbf{x}) = \mathbf{y}_{jk} - \mathbf{s}(\mathbf{T}_k \mathbf{p}_j) \quad (179)$$

- all the errors can be stacked up into a tall **error vector** as

$$\mathbf{e}_y(\mathbf{x}_{\text{op}}) = \begin{bmatrix} \mathbf{e}_{y,10}(\mathbf{x}_{\text{op}}) \\ \mathbf{e}_{y,20}(\mathbf{x}_{\text{op}}) \\ \vdots \\ \mathbf{e}_{y,MK}(\mathbf{x}_{\text{op}}) \end{bmatrix} \quad (180)$$

Linearizing \mathbf{T}_p

- it is worth a small aside to see how to **linearize the product \mathbf{T}_p**
- we insert our various **perturbation schemes** so that

$$\begin{aligned}\mathbf{T}_p &= \exp(\epsilon^\wedge) \mathbf{T}_{op} (\mathbf{p}_{op} + \mathbf{D}\zeta) \\ &\approx (\mathbf{1} + \epsilon^\wedge) \mathbf{T}_{op} (\mathbf{p}_{op} + \mathbf{D}\zeta) \\ &\approx \mathbf{T}_{op}\mathbf{p}_{op} + \epsilon^\wedge \mathbf{T}_{op}\mathbf{p}_{op} + \mathbf{T}_{op}\mathbf{D}\zeta + \underbrace{\epsilon^\wedge \mathbf{T}_{op}\mathbf{D}\zeta}_{\text{very small}} \\ &\approx \mathbf{T}_{op}\mathbf{p}_{op} + (\mathbf{T}_{op}\mathbf{p}_{op})^\odot \epsilon + \mathbf{T}_{op}\mathbf{D}\zeta \\ &= \mathbf{T}_{op}\mathbf{p}_{op} + [(\mathbf{T}_{op}\mathbf{p}_{op})^\odot \quad \mathbf{0} \quad \mathbf{T}_{op}\mathbf{D}] \begin{bmatrix} \epsilon \\ \eta \\ \zeta \end{bmatrix}\end{aligned}$$

- then we can use the **chain rule** to write

$$\mathbf{s}(\mathbf{T}_p) \approx \mathbf{s}(\mathbf{T}_{op}\mathbf{p}_{op}) + \mathbf{S} [(\mathbf{T}_{op}\mathbf{p}_{op})^\odot \quad \mathbf{0} \quad \mathbf{T}_{op}\mathbf{D}] \begin{bmatrix} \epsilon \\ \eta \\ \zeta \end{bmatrix} \quad (182)$$

with \mathbf{S} the Jacobian of the sensor model, $\mathbf{s}(\cdot)$

MAP cost

- **MAP cost** is then $J(\mathbf{x}) = J_v(\mathbf{x}) + J_y(\mathbf{x})$, which can be quadratically approximated as

$$J(\mathbf{x}) = J_v(\mathbf{x}) + J_y(\mathbf{x}) \approx J(\mathbf{x}_{\text{op}}) - \mathbf{e}(\mathbf{x}_{\text{op}})^T \mathbf{W}^{-1} \mathbf{H} \delta \mathbf{x} + \delta \mathbf{x}^T \mathbf{H}^T \mathbf{W}^{-1} \mathbf{H} \delta \mathbf{x} \quad (185)$$

with

$$\mathbf{e}(\mathbf{x}_{\text{op}}) = \begin{bmatrix} \mathbf{e}_v(\mathbf{x}_{\text{op}}) \\ \mathbf{e}_y(\mathbf{x}_{\text{op}}) \end{bmatrix}, \quad \delta \mathbf{x} = \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{G}_1 & \mathbf{G}_2 \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix}$$

- the **minimizing perturbation**, $\delta \mathbf{x}^*$, is the solution to

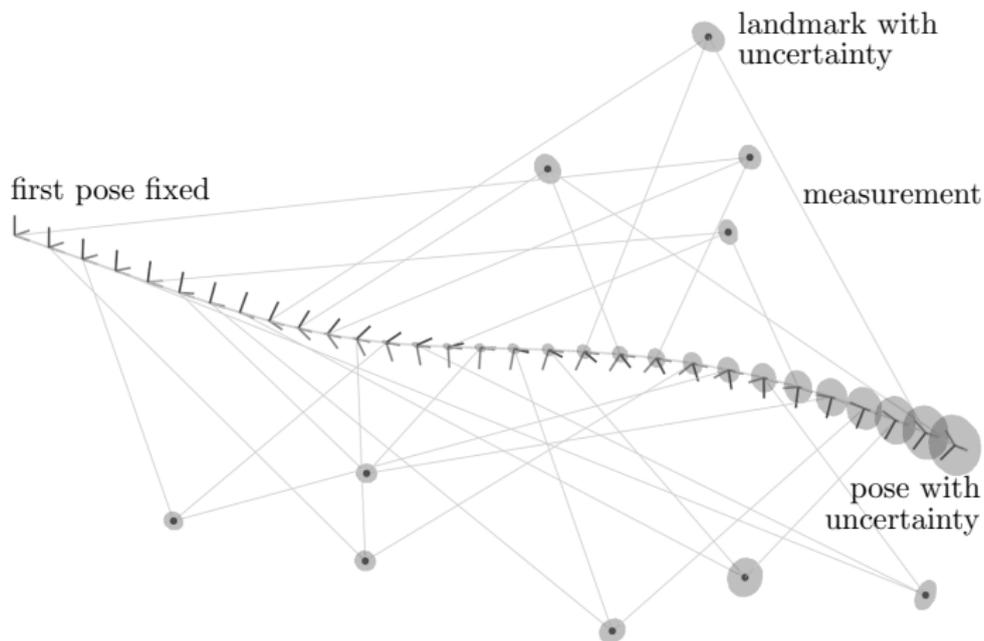
$$\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H} \delta \mathbf{x}^* = \mathbf{H}^T \mathbf{W}^{-1} \mathbf{e}(\mathbf{x}_{\text{op}}) \quad (186)$$

- we solve for $\delta \mathbf{x}^*$, then apply the optimal perturbations

$$\mathbf{T}_{\text{op},k} \leftarrow \exp(\hat{\boldsymbol{\epsilon}}_k^*) \mathbf{T}_{\text{op},k}, \quad \boldsymbol{\varpi}_{\text{op},k} \leftarrow \boldsymbol{\varpi}_{\text{op},k} + \boldsymbol{\eta}_k^*, \quad \mathbf{p}_{\text{op},j} \leftarrow \mathbf{p}_{\text{op},j} + \mathbf{D} \boldsymbol{\zeta}_j^*$$

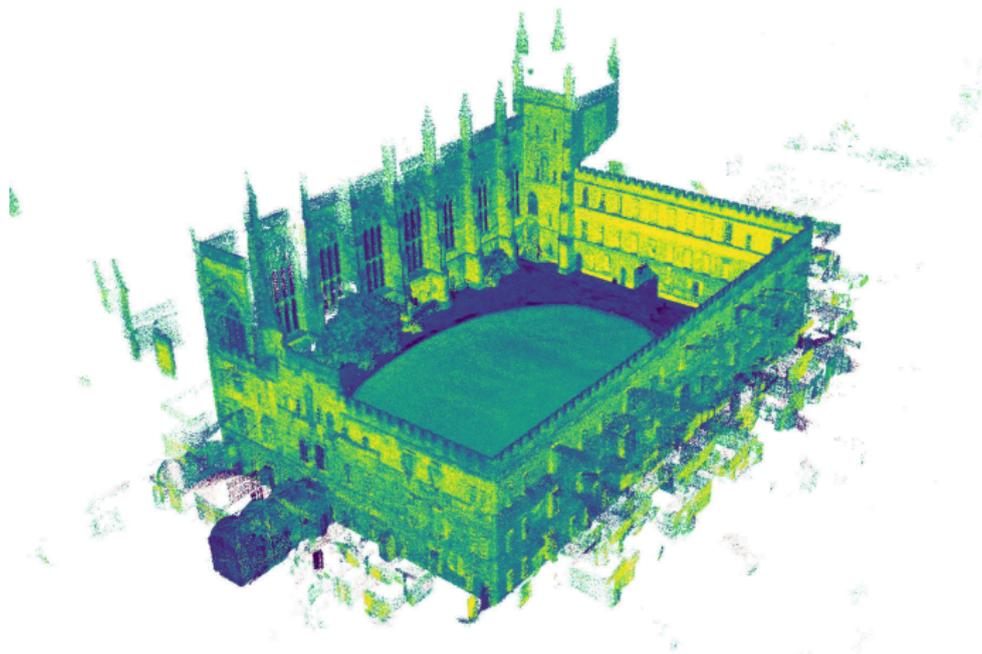
and iterate to convergence

STEAM example



- with only **one landmark measurement per pose**, the motion prior is necessary to solve

STEAM example



- these continuous-time methods can be used as a **sliding-window filter** and run in realtime

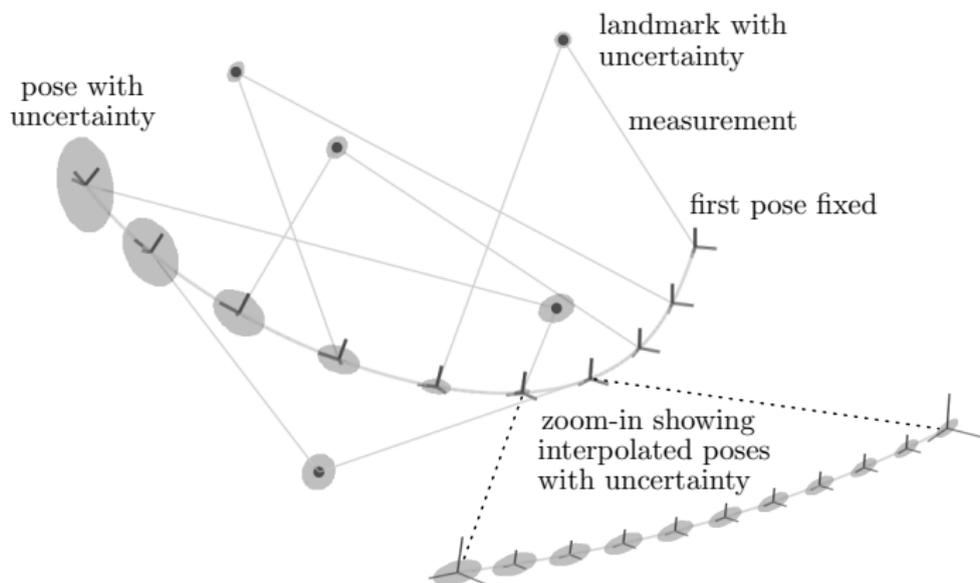
STEAM sparsity

- the left-hand side in our linear system of equations has the form

$$\mathbf{H}^T \mathbf{W}^{-1} \mathbf{H} = \begin{bmatrix} \mathbf{F}^{-T} \mathbf{Q}^{-1} \mathbf{F}^{-1} + \mathbf{G}_1^T \mathbf{R}^{-1} \mathbf{G}_1 & \mathbf{G}_1^T \mathbf{R}^{-1} \mathbf{G}_2 \\ \mathbf{G}_2^T \mathbf{R}^{-1} \mathbf{G}_1 & \mathbf{G}_2^T \mathbf{R}^{-1} \mathbf{G}_2 \end{bmatrix} \quad (187)$$

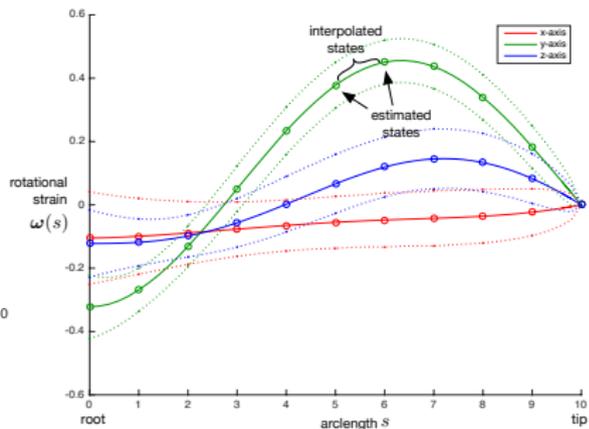
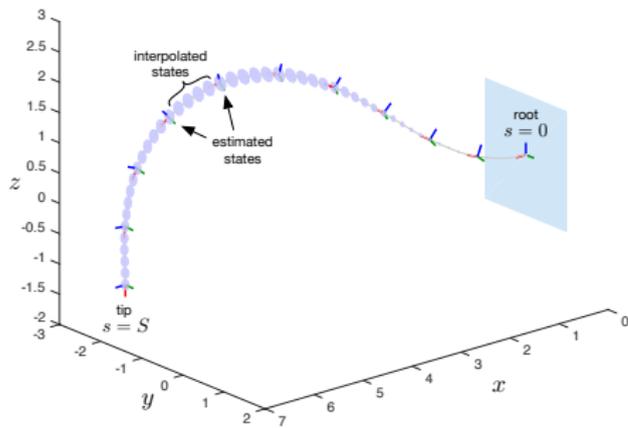
- this preserves the usual **SLAM arrowhead form**
 - the top-left block is **block-tridiagonal**
 - the bottom-right block is **block-diagonal**
- a **sparse Cholesky decomposition** or **Schur complement** can be used to solve the linear system efficiently

STEAM interpolation



- we can still use **GP interpolation** after the main solve – see book for details

Continuum robot example



- we can use these tools to estimate the shape of **snake-like** robots
- we simply replace time, t , with **arclength**, s

Lecture Summary

- we can build useful motion priors starting from **physically motivated** stochastic differential equations
- we then reframe estimation as **Gaussian process regression**
- still very efficient due to our particular motion priors / kernels – $O(K)$ for main solve, $O(1)$ for a query – reason is the **Markov property**
- ideas can be applied even when the state variables are members of **Lie groups**
- ideas can be used within a SLAM setup – **STEAM**
- **Lecture 4:** Beyond MAP – Variational Inference

Lecture 4: Beyond MAP – Variational Inference

A Short Course in State Estimation

Timothy D. Barfoot

University of Toronto

Copyright © 2024

Lecture Outline

Lecture 1: Linear-Gaussian Estimation

Lecture 2: Nonlinearities and Lie Groups

Lecture 3: Continuous-Time Estimation

Lecture 4: Beyond MAP – Variational Inference

- Motivation

- Variational Inference

- Sparsity

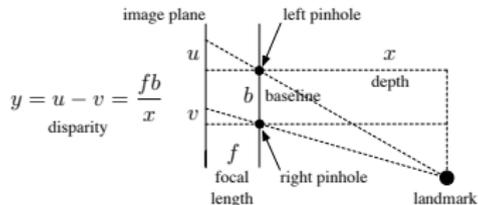
- Derivative-free

- Parameter Learning

Motivation

- we have primarily presented an **MAP** approach to state estimation, and this certainly has its place
- but, is it possible to do better than that? can we be better **Bayesians**?
- consider a one-dimensional example where we get a single **disparity** measurement from a **stereo camera**:

$$y = \frac{fb}{x} + n, \quad n \sim \mathcal{N}(0, R)$$



(188)

- if we have a prior on position, $p(x) = \mathcal{N}(\check{x}, \check{P})$, what is the **Bayesian posterior** once we incorporate this measurement?

Stereo camera example

- to perform **Bayesian inference**,

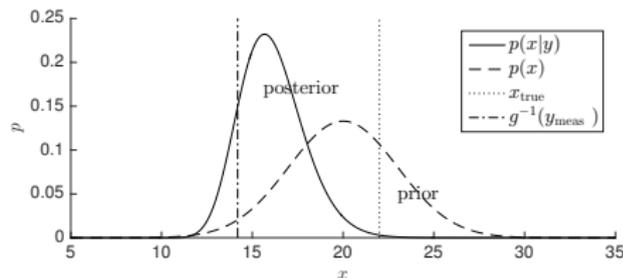
$$p(x|y) = \frac{p(y|x)p(x)}{\int_{-\infty}^{\infty} p(y|x)p(x) dx} \quad (189)$$

we require expressions for $p(y|x)$ and $p(x)$

$$p(y|x) = \mathcal{N}\left(\frac{fb}{x}, R\right) = \frac{1}{\sqrt{2\pi R}} \exp\left(-\frac{1}{2R}\left(y - \frac{fb}{x}\right)^2\right), \quad p(x) = \mathcal{N}(\tilde{x}, \tilde{P}) = \frac{1}{\sqrt{2\pi\tilde{P}}} \exp\left(-\frac{1}{2\tilde{P}}(x - \tilde{x})^2\right)$$

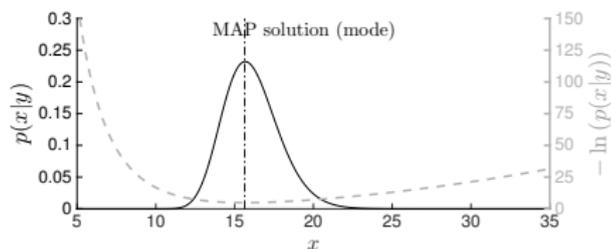
$$\tilde{x} = 20 \text{ [m]}, \quad \tilde{P} = 9 \text{ [m}^2\text{]}, \quad f = 400 \text{ [pixel]}, \quad b = 0.1 \text{ [m]}, \quad R = 0.09 \text{ [pixel}^2\text{]}, \quad x_{\text{true}} = 22 \text{ [m]}, \quad y_{\text{meas}} = \frac{fb}{x_{\text{true}}} + 1 \text{ [pixel]}$$

- the **Bayesian posterior** then looks like this (computed numerically)

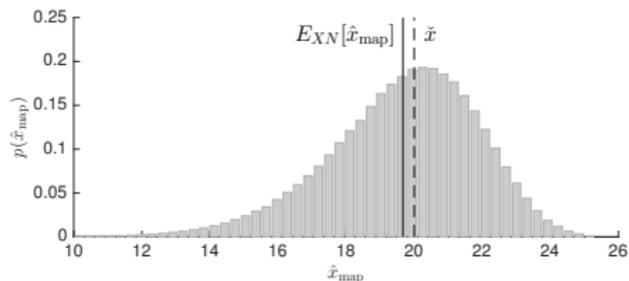


Stereo camera example

- clearly the Bayesian posterior is **not Gaussian**
- the **MAP** solution will find the **mode** of the Bayesian posterior



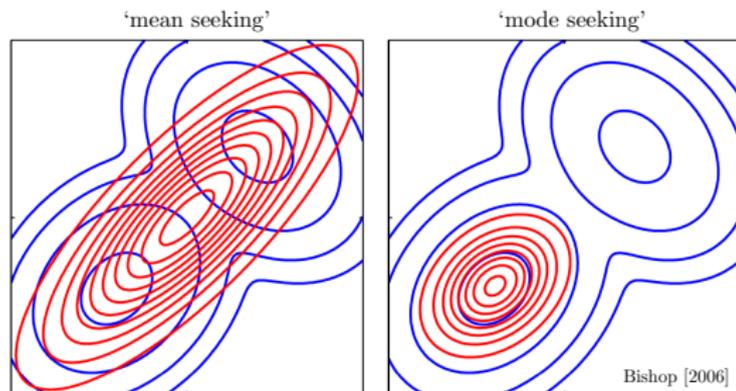
- is this what we want? we end up with a -33.0 cm average **bias**



- we might prefer an algorithm that can estimate the **mean** of the posterior if this is the metric we care about

Variational inference

- we could try **variational (Bayesian) inference**
- find a (Gaussian) estimate, $q(\mathbf{x})$, that is 'closest' to the Bayesian posterior, $p(\mathbf{x}|\mathbf{z})$, in terms of the **Kullback-Leibler divergence**



$$\begin{aligned} \text{KL}(p||q) &= - \int_{-\infty}^{\infty} p(\mathbf{x}|\mathbf{z}) \ln \left(\frac{q(\mathbf{x})}{p(\mathbf{x}|\mathbf{z})} \right) d\mathbf{x} \\ &= E_p [\ln p(\mathbf{x}|\mathbf{z}) - \ln q(\mathbf{x})] \end{aligned}$$

$$\begin{aligned} \text{KL}(q||p) &= - \int_{-\infty}^{\infty} q(\mathbf{x}) \ln \left(\frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{x})} \right) d\mathbf{x} \\ &= E_q [\ln q(\mathbf{x}) - \ln p(\mathbf{x}|\mathbf{z})], \end{aligned}$$

- we pick the KL divergence on the right since easier to compute

Gaussian variational inference

- we assume our estimate is **Gaussian**

$$q(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (190)$$

- our chosen **KL divergence** then simplifies to

$$\text{KL}(q||p) = E_q[-\ln p(\mathbf{x}, \mathbf{z})] - \underbrace{\frac{1}{2} \ln((2\pi e)^N |\boldsymbol{\Sigma}|)}_{\text{entropy}} + \underbrace{\ln p(\mathbf{z})}_{\text{constant}} \quad (191)$$

- dropping the constant term and defining $\phi(\mathbf{x}) = -\ln p(\mathbf{x}, \mathbf{z})$ we can define a **functional** as

$$V(q) = E_q[\phi(\mathbf{x})] + \frac{1}{2} \ln(|\boldsymbol{\Sigma}^{-1}|) \quad (192)$$

Gaussian variational inference

- our **functional** is

$$V(q) = E_q[\phi(\mathbf{x})] + \frac{1}{2} \ln (|\Sigma^{-1}|) \quad (193)$$

- the derivatives with respect to the unknown mean, $\boldsymbol{\mu}$, and inverse covariance, Σ^{-1} , are

$$\frac{\partial V(q)}{\partial \boldsymbol{\mu}^T} = \Sigma^{-1} E_q[(\mathbf{x} - \boldsymbol{\mu})\phi(\mathbf{x})]$$

$$\frac{\partial^2 V(q)}{\partial \boldsymbol{\mu}^T \partial \boldsymbol{\mu}} = \Sigma^{-1} E_q[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \phi(\mathbf{x})] \Sigma^{-1} - \Sigma^{-1} E_q[\phi(\mathbf{x})]$$

$$\frac{\partial V(q)}{\partial \Sigma^{-1}} = -\frac{1}{2} E_q[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \phi(\mathbf{x})] + \frac{1}{2} \Sigma E_q[\phi(\mathbf{x})] + \frac{1}{2} \Sigma$$

Gaussian variational inference

- with these derivatives we can **approximate** the functional as

$$V(q^{(i+1)}) \approx V(q^{(i)}) + \left(\frac{\partial V(q)}{\partial \boldsymbol{\mu}^T} \Big|_{q^{(i)}} \right)^T \delta \boldsymbol{\mu} + \frac{1}{2} \delta \boldsymbol{\mu}^T \left(\frac{\partial^2 V(q)}{\partial \boldsymbol{\mu}^T \partial \boldsymbol{\mu}} \Big|_{q^{(i)}} \right) \delta \boldsymbol{\mu} + \text{tr} \left(\frac{\partial V(q)}{\partial \boldsymbol{\Sigma}^{-1}} \Big|_{q^{(i)}} \delta \boldsymbol{\Sigma}^{-1} \right) \quad (195)$$

which we now want to iteratively minimize with respect to

$$\delta \boldsymbol{\mu} = \boldsymbol{\mu}^{(i+1)} - \boldsymbol{\mu}^{(i)} \text{ and } \delta \boldsymbol{\Sigma}^{-1} = (\boldsymbol{\Sigma}^{-1})^{(i+1)} - (\boldsymbol{\Sigma}^{-1})^{(i)}$$

- clearly we want $\delta \boldsymbol{\mu}$ to be the solution to

$$\underbrace{\left(\frac{\partial^2 V(q)}{\partial \boldsymbol{\mu}^T \partial \boldsymbol{\mu}} \Big|_{q^{(i)}} \right)}_{(\boldsymbol{\Sigma}^{-1})^{(i+1)}} \delta \boldsymbol{\mu} = - \left(\frac{\partial V(q)}{\partial \boldsymbol{\mu}^T} \Big|_{q^{(i)}} \right) \quad (196)$$

- the indicated inverse-covariance update results from setting $\frac{\partial V(q)}{\partial \boldsymbol{\Sigma}^{-1}}$ to zero and noticing the equivalence

$$\frac{\partial^2 V(q)}{\partial \boldsymbol{\mu}^T \partial \boldsymbol{\mu}} = \boldsymbol{\Sigma}^{-1} - 2 \boldsymbol{\Sigma}^{-1} \frac{\partial V(q)}{\partial \boldsymbol{\Sigma}^{-1}} \boldsymbol{\Sigma}^{-1} \quad (197)$$

Natural gradient descent interpretation

- it turns out this update scheme is equivalent to **natural gradient descent**

$$\delta\alpha = -\mathcal{I}_\alpha^{-1} \frac{\partial V(q)}{\partial \alpha^T} \quad (198)$$

where

$$\alpha = \begin{bmatrix} \boldsymbol{\mu} \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}, \quad \delta\alpha = \begin{bmatrix} \delta\boldsymbol{\mu} \\ \text{vec}(\delta\boldsymbol{\Sigma}^{-1}) \end{bmatrix}, \quad \frac{\partial V(q)}{\partial \alpha^T} = \begin{bmatrix} \frac{\partial V(q)}{\partial \boldsymbol{\mu}^T} \\ \text{vec} \left(\frac{\partial V(q)}{\partial \boldsymbol{\Sigma}^{-1}} \right) \end{bmatrix}, \quad \mathcal{I}_\alpha = \begin{bmatrix} \boldsymbol{\Sigma}^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) \end{bmatrix}$$

- importantly, \mathcal{I}_α is the **Fisher information matrix** for the variational parameter, α
- NGD is essentially a **quasi-Newton method** that uses the FIM as an approximation of the Hessian

Stein's lemma to the rescue

- we can then employ **Stein's lemma**

$$E_q[(\mathbf{x} - \boldsymbol{\mu})\phi(\mathbf{x})] = \boldsymbol{\Sigma} E_q \left[\frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}^T} \right] \quad (199)$$

to clean up the update further so that it becomes

$$(\boldsymbol{\Sigma}^{-1})^{(i+1)} = E_{q^{(i)}} \left[\frac{\partial^2}{\partial \mathbf{x}^T \partial \mathbf{x}} \phi(\mathbf{x}) \right] \quad (200a)$$

$$(\boldsymbol{\Sigma}^{-1})^{(i+1)} \delta \boldsymbol{\mu} = -E_{q^{(i)}} \left[\frac{\partial}{\partial \mathbf{x}^T} \phi(\mathbf{x}) \right] \quad (200b)$$

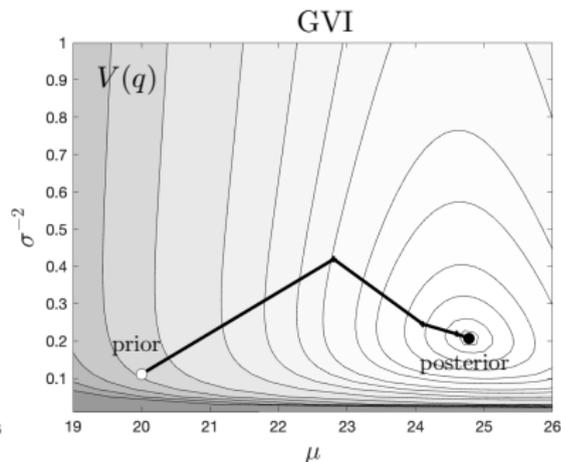
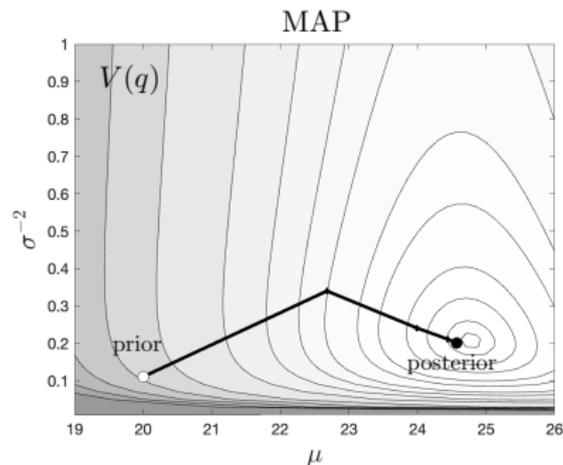
$$\boldsymbol{\mu}^{(i+1)} = \boldsymbol{\mu}^{(i)} + \delta \boldsymbol{\mu} \quad (200c)$$

where i is the iteration index

- this looks just like an **MAP** update via **Newton's method** – but now we have **expectations** wrapping the derivatives of $\phi(\mathbf{x})$

Stereo camera example

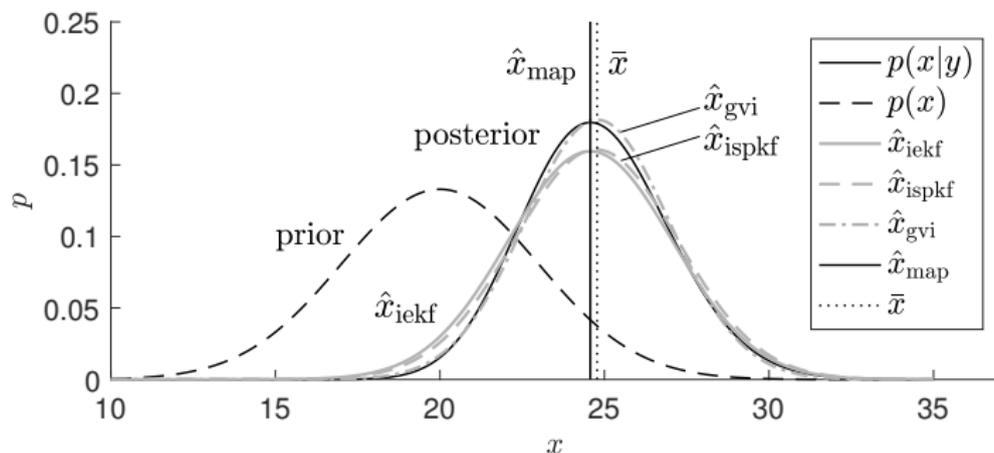
- if we return to our **stereo camera example**, we can compare the performance of **MAP** and **GVI**



- **MAP** approximates the expectations, $E_q[\cdot]$, only at the **mean** and as a result does not get to the minimum of $V(q)$

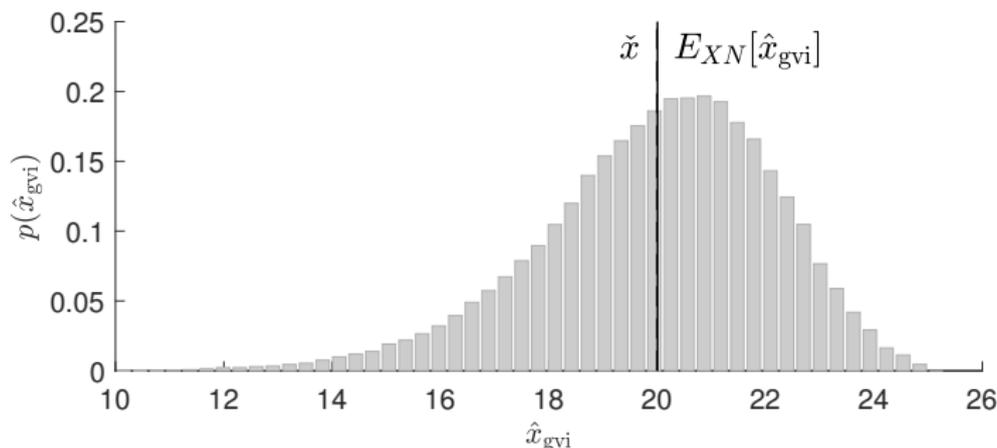
Stereo camera example

- we can see that **GVI** also converges to a point where its mean looks closer to the **mean of the true posterior**



Stereo camera example

- histogram of the errors over many trials shows **GVI** has an average error of 0.28 cm whereas **MAP** was -33.0 cm



- **GVI** seems to help remove the bias in this controlled example

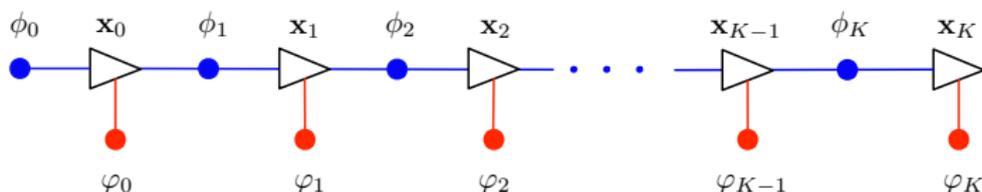
Can we scale up?

- there are some expectations in our update scheme that look **computationally expensive**

$$\underbrace{E_q[\phi(\mathbf{x})]}_{\text{scalar}}, \quad \underbrace{E_q\left[\frac{\partial}{\partial \mathbf{x}^T} \phi(\mathbf{x})\right]}_{\text{column}}, \quad \underbrace{E_q\left[\frac{\partial^2}{\partial \mathbf{x}^T \partial \mathbf{x}} \phi(\mathbf{x})\right]}_{\text{matrix}} \quad (201)$$

- however, if we assume there is some sparsity (in the usual **factor graph** way) then our negative log-likelihood decomposes as

$$\phi(\mathbf{x}) = \sum_{k=1}^K \phi_k(\mathbf{x}_k) \quad (202)$$



Can we scale up?

- our negative log-likelihood decomposes as

$$\phi(\mathbf{x}) = \sum_{k=1}^K \phi_k(\mathbf{x}_k) \quad (203)$$

- inserting this we see for the first expectation that

$$E_q[\phi(\mathbf{x})] = E_q \left[\sum_{k=1}^K \phi_k(\mathbf{x}_k) \right] = \sum_{k=1}^K E_q[\phi_k(\mathbf{x}_k)] = \sum_{k=1}^K E_{q_k}[\phi_k(\mathbf{x}_k)] \quad (204)$$

where the expectations are now only over $q_k(\mathbf{x}_k) = \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_{kk})$, the **marginal** for the variables involved in $\phi_k(\mathbf{x}_k)$

- similar simplifications occur for the other two expectations

$$E_q \left[\frac{\partial}{\partial \mathbf{x}^T} \phi(\mathbf{x}) \right] = \sum_{k=1}^K \mathbf{P}_k^T E_{q_k} \left[\frac{\partial}{\partial \mathbf{x}_k^T} \phi_k(\mathbf{x}_k) \right] \quad (205a)$$

$$E_q \left[\frac{\partial^2}{\partial \mathbf{x}^T \partial \mathbf{x}} \phi(\mathbf{x}) \right] = \sum_{k=1}^K \mathbf{P}_k^T E_{q_k} \left[\frac{\partial^2}{\partial \mathbf{x}_k^T \partial \mathbf{x}_k} \phi_k(\mathbf{x}_k) \right] \mathbf{P}_k \quad (205b)$$

where \mathbf{P}_k is a **selection matrix** so that $\mathbf{x}_k = \mathbf{P}_k \mathbf{x}$

Can we scale up?

- with this simplification of the expectations, we call our scheme **exactly sparse Gaussian variational inference (ESGVI)** – we only need to compute expectations over the **marginals** associated with the variables in each factor
- this means that we only need the blocks of the covariance, Σ , associated with the **non-zero blocks** of the inverse covariance, Σ^{-1}
- Takahashi et al. [1973] shows how to compute a set of covariance entries that is just what we need, plus some extra blocks that come from **fill-in** – e.g., Cholesky decomposition $\Sigma^{-1} = \mathbf{LL}^T$
- finding a variable order that minimizes fill-in is known to be an NP-hard problem in general
- we piggyback finding the required covariance entries onto the usual **MAP** solve, so **ESGVI** has the same big-O complexity as **MAP**, but the expectations still make it slower

Fill-in

- typical **fill-in** '+' in for our estimation problems – we efficiently compute blocks of Σ where \mathbf{L} is non-zero

basic sparsity constraint
(note fill in at (5,3) in \mathbf{L})

$$\Sigma^{-1} = \begin{bmatrix} * & & * & & * \\ & * & & & \\ * & & * & & \\ & & & * & \\ * & & & & * \\ & & & & & * \end{bmatrix}$$

$$\mathbf{L} = \begin{bmatrix} * & & & & & \\ & * & & & & \\ * & & * & & & \\ & & & * & & \\ * & & & & * & \\ & & & & & * \end{bmatrix}$$

trajectory example
(6 robot poses)

$$\Sigma^{-1} = \begin{bmatrix} * & * & & & & \\ * & * & * & & & \\ & * & * & * & & \\ & & * & * & * & \\ & & & * & * & * \\ & & & & * & * \end{bmatrix}$$

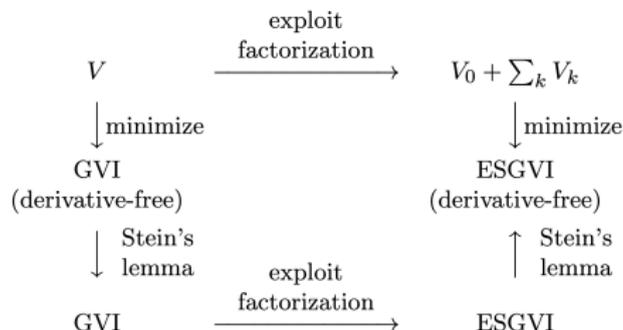
$$\mathbf{L} = \begin{bmatrix} * & & & & & \\ * & * & & & & \\ & * & * & & & \\ & & * & * & & \\ & & & * & * & \\ & & & & * & * \end{bmatrix}$$

SLAM example
(3 poses, 3 landmarks)

$$\Sigma^{-1} = \begin{bmatrix} * & * & & * & * & * \\ * & * & * & * & * & * \\ & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix}$$

$$\mathbf{L} = \begin{bmatrix} * & & & & & \\ * & * & & & & \\ & * & * & & & \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix}$$

Derivative-free version



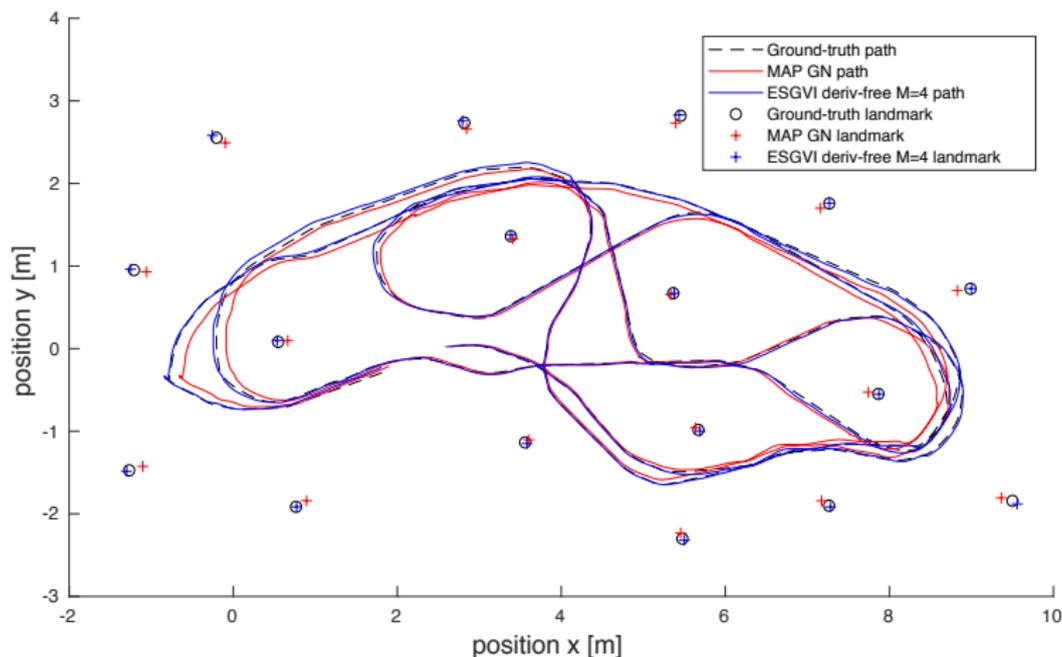
- we can also use **Stein's lemma** (opposite direction) at the individual factor level to compute all the small expectations **derivative-free**

$$\text{e.g., } E_{q_k} \left[\frac{\partial}{\partial \mathbf{x}_k^T} \phi_k(\mathbf{x}_k) \right] = \Sigma_{kk}^{-1} E_{q_k} [(\mathbf{x}_k - \boldsymbol{\mu}_k) \phi_k(\mathbf{x}_k)] \quad (206)$$

- this means we can do **batch state estimation without derivatives!**

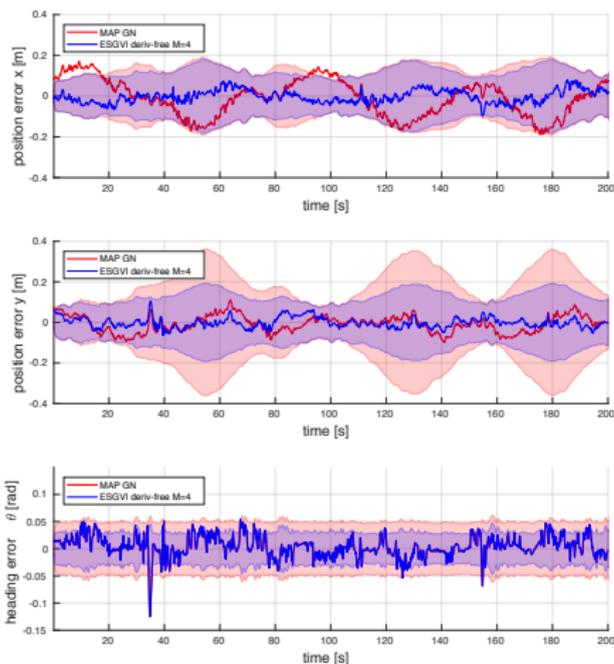
Example: bearing-only SLAM

- we can find challenging estimation cases where **GVI** outperforms **MAP** – the gains are usually small



Example: bearing-only SLAM

- we can find challenging estimation cases where **GVI** outperforms **MAP** – the gains are usually small but measurable



Parameter learning

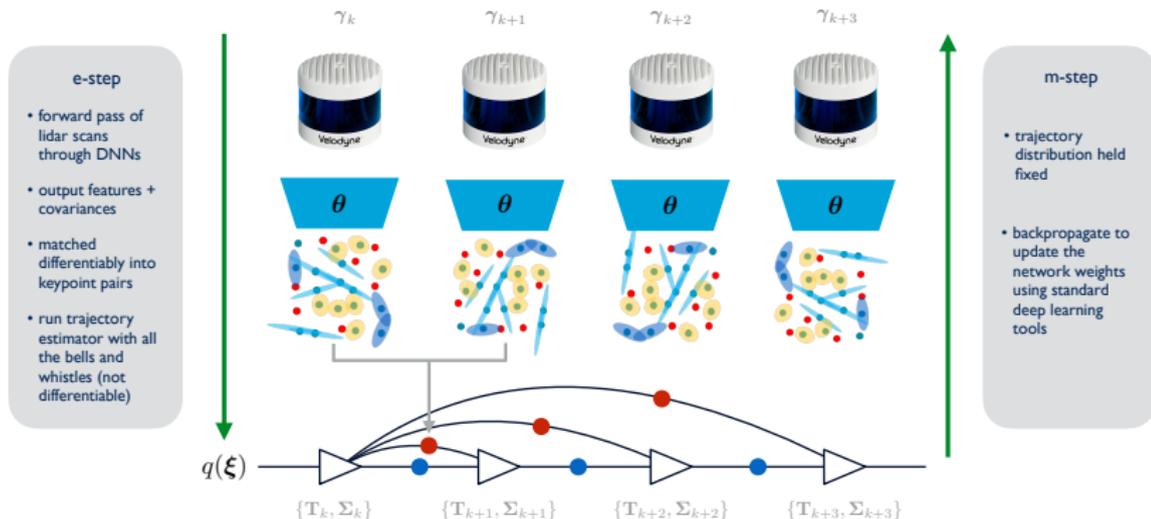
- an added bonus of this framework is that we can also **learn parameters, θ**
- the negative log-likelihood of our data (given the parameters) can be written as

$$-\ln(p(\mathbf{z}|\theta)) = \underbrace{\int_{-\infty}^{\infty} q(\mathbf{x}) \ln \left(\frac{p(\mathbf{x}|\mathbf{z}, \theta)}{q(\mathbf{x})} \right) d\mathbf{x}}_{-\text{KL}(q||p) \leq 0} - \underbrace{\int_{-\infty}^{\infty} q(\mathbf{x}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{x})} \right) d\mathbf{x}}_{V(q|\theta)} \quad (207)$$

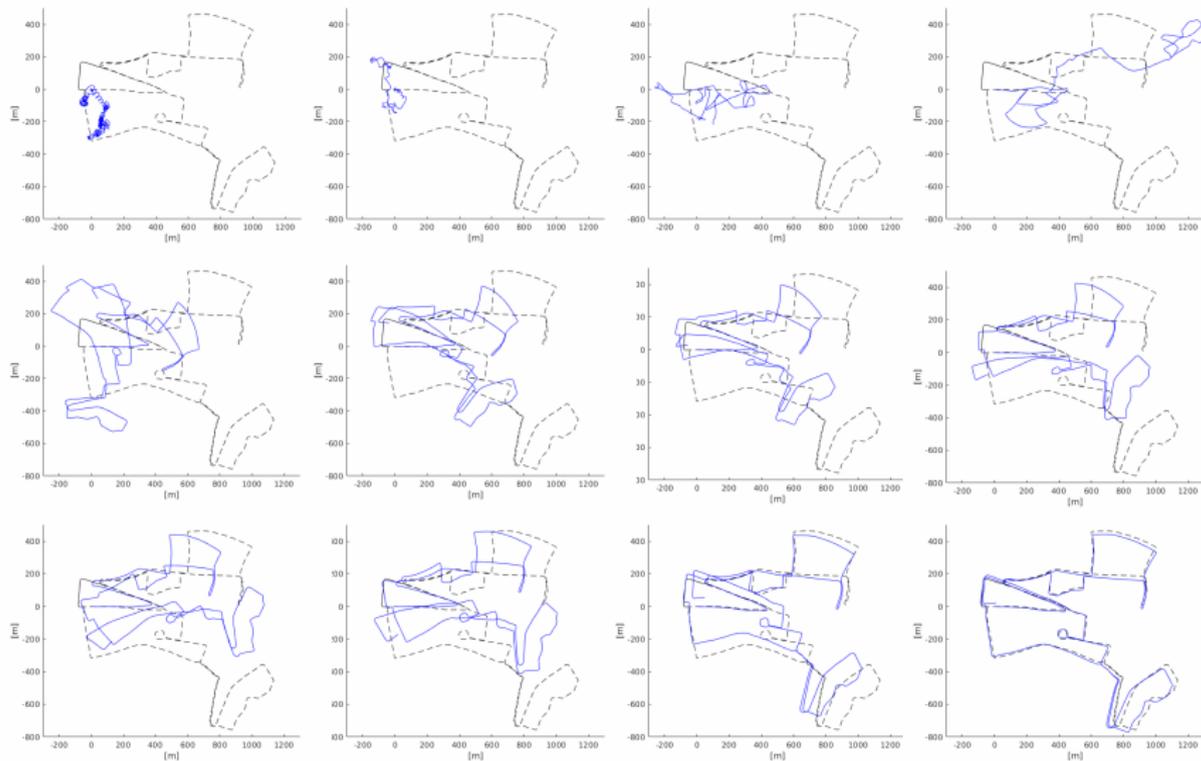
- we can use **expectation minimization** [Neal and Hinton, 1998]:
 - **e-step**: hold θ fixed, minimize $V(q|\theta)$ w.r.t. q
 - **m-step**: hold q fixed, minimize $V(q|\theta)$ w.r.t. θ
- parameters can be anything – measurement covariances, calibrations, network weights for measurement model

Example: parameter learning

- we can do **unsupervised** parameter learning with our full **non-differentiable** state estimator

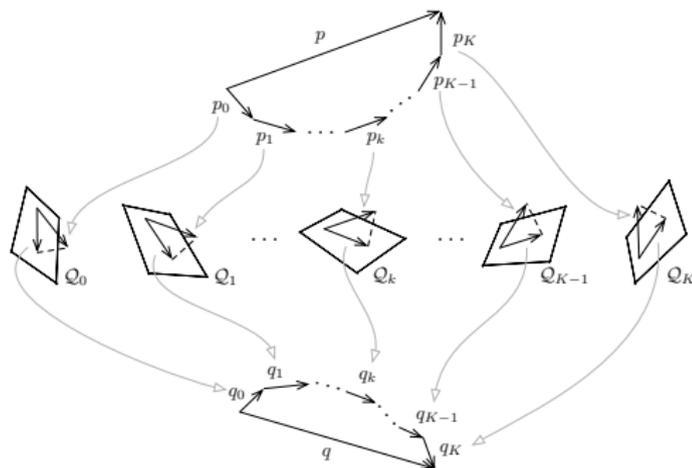


Example: parameter learning



Geometric interpretation

- can also view **ESGVI** through the lens of **Bayes space** – a high-dimensional geometric space where entire PDFs can be viewed as a single **vector**



- inference is then **iterative projection** of the full Bayesian posterior onto a subspace containing Gaussians (Barfoot and D'Eleuterio, 2023)

Lecture Summary

- we have seen that **Gaussian variational inference** can be used to produce a ‘better’ solution than **MAP**
- we can still exploit **sparsity** to make **ESGVI** have the same big-O complexity as **MAP** – but compute time is still a large factor slower
- we can avoid even computing derivatives by using **Stein’s lemma** and **Gaussian cubature** to evaluate expectations
- in practice we can run **MAP** first since it’s faster, then refine the answer with **ESGVI** for a few iterations
- can fold in **parameter learning** via **expectation minimization**
- there is a way to also enforce that Σ remains positive definite by exploiting the connection to **NGD** [Goudar et al., 2022]
- although not discussed, can make the ideas work on **Lie groups**
- **ESGVI** is related to but distinct from **Stein variational inference**

Course Conclusion

- looking for next steps?
- estimation, solvers:
 - general overview (Thrun et al., 2006; Barfoot, 2024)
 - **incremental smoothing and mapping (iSAM2)** (Kaess et al., 2012)
 - **certifiably optimal estimation** (Yang and Carlone, 2022)
- Lie groups:
 - **invariant EKF** (Bonnabel et al., 2008)
 - **equivariant EKF** (Mahony and Trumppf, 2021)
- variational inference:
 - **Stein variational inference** (Liu and Wang, 2016; Maken et al., 2022)

References

- Barfoot, T. D., *State Estimation for Robotics*, Cambridge University Press, 2nd edition, 2024.
- Barfoot, T. D. and D'Eleuterio, G. M. T., "Variational Inference as Iterative Projection in a Bayesian Hilbert Space with Application to Robotic State Estimation," *Robotica*, 41(2):632–667, 2023.
- Bishop, C. M., *Pattern Recognition and Machine Learning*, Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- Bonnabel, S., Martin, P., and Rouchon, P., "Symmetry-preserving observers," *IEEE Transactions on Automatic Control*, 53(11):2514–2526, 2008.
- Boumal, N., *An introduction to optimization on smooth manifolds*, Cambridge University Press, 2023.
- Goudar, A., Zhao, W., Barfoot, T. D., and Schoellig, A. P., "Gaussian Variational Inference with Covariance Constraints Applied to Range-Only Localization," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robot Systems (IROS)*, Kyoto, Japan, 2022.
- Kaess, M., Johannsson, H., Roberts, R., Ila, V., Leonard, J. J., and Dellaert, F., "iSAM2: Incremental Smoothing and Mapping Using the Bayes Tree," *IJRR*, 31(2):217–236, 2012.
- Liu, Q. and Wang, D., "Stein variational gradient descent: A general purpose bayesian inference algorithm," *Advances in neural information processing systems*, 29, 2016.
- Mahony, R. and Trumpf, J., "Equivariant filter design for kinematic systems on lie groups," *IFAC-PapersOnLine*, 54(9):253–260, 2021.
- Maken, F. A., Ramos, F., and Ott, L., "Stein particle filter for nonlinear, non-Gaussian state estimation," *IEEE Robotics and Automation Letters*, 7(2):5421–5428, 2022.
- Rasmussen, C. E. and Williams, C. K. I., *Gaussian Processes for Machine Learning*, MIT Press, Cambridge, MA, 2006.
- Takahashi, K., Fagan, J., and Chen, M.-S., "A Sparse Bus Impedance Matrix and its Application to Short Circuit Study," in *Proceedings of the PICA Conference*, 1973.
- Thrun, S., Burgard, W., and Fox, D., *Probabilistic Robotics*, MIT Press, 2006.
- Yang, H. and Carlone, L., "Certifiably optimal outlier-robust geometric perception: Semidefinite relaxations and scalable global optimization," *IEEE transactions on pattern analysis and machine intelligence*, 45(3):2816–2834, 2022.