

# Random Features, Density Estimate and Simultaneous Localization and Mapping

Yermek Kapushev

y.kapushev@skoltech.ru

Skolkovo Institute of Science and Technology  
Moscow, Russia

Moscow, 2020

Random Features

Denoising Score Matching

Simultaneous Localization and Mapping

Summary

► Data set

$$(\mathbf{X}, \mathbf{y}) = \{(x_1, y_1), \dots, (x_n, y_n)\} \in \mathcal{X} \times \mathbb{R}, \quad y_i = f(x_i) + \varepsilon_i, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

► Let  $\phi(x) : \mathcal{X} \rightarrow \mathbb{R}^d$  be some feature map

► Suppose that  $y = \beta^\top \phi(x)$ , then prediction is given by

$$\hat{f}(x^*) = \phi(x^*)^\top \left( \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda \mathbf{I} \right)^{-1} \phi(\mathbf{X})^\top \mathbf{y}, \quad \mathcal{O}(d^3)$$

► Let  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  be a *kernel function*. Then prediction can be rewritten as

$$\hat{f}(x^*) = \mathbf{k}_*^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}, \quad \mathcal{O}(n^3)$$

where  $\mathbf{k}_* = (k(x^*, x_1), \dots, k(x^*, x_n))$ ,  $\mathbf{K}_{ij} = k(x_i, x_j)$ .

# Random Fourier Features (RFF)

## Theorem (Bochner)

*A continuous kernel  $k(x, x') = k(x - x')$  on  $\mathbb{R}^d$  is positive definite if and only if  $k(\delta)$  is a Fourier transform of a non-negative measure*

$$k(x, x') = \int_{\Omega} p(w) e^{jw^{\top}(x-x')} dw$$

**Idea:** use Monte-Carlo to approximate the integral<sup>1</sup>

$$k_{RFF}(x, x') = \frac{1}{D} \sum_{i=1}^D \cos(w_i^\top (x - x')) = \psi_{\mathbf{W}}(x)^\top \psi_{\mathbf{W}}(x'),$$

where

- $w_i \sim p(w)$
- $\mathbf{W} = (w_1^\top, \dots, w_D^\top)^\top$
- $\psi_{\mathbf{W}}(x) = 1/\sqrt{D}(\cos(w_1^\top x), \sin(w_1^\top x), \dots, \cos(w_D^\top x), \sin(w_D^\top x))^\top$

$$\hat{\mathbf{K}} = \Psi \Psi^\top, \quad \Psi = \|\psi_{\mathbf{W}}(x_i)^\top\|_{i=1}^n$$

→ Go back to linear model with  $\psi_{\mathbf{W}}(x)$  features

---

<sup>1</sup>Rahimi, A., Recht, B. (2008). Random features for large-scale kernel machines.

- ▶ Quasi Monte-Carlo (QMC)<sup>2</sup>
- ▶ Gaussian Quadrature <sup>3</sup>
- ▶ Orthogonal Random Features (ORF)<sup>4</sup>
- ▶ Random Orthogonal Matrices (ROM)<sup>5</sup>
- ▶ Ridge Leverage Score based Features <sup>6</sup>

---

<sup>2</sup>Yang, et al. (2016). Quasi-Monte Carlo feature maps for shift-invariant kernels.

<sup>3</sup>Dao et al. (2017). Gaussian quadrature for kernel features.

<sup>4</sup>Felix, X. Yu, et al. (2016). Orthogonal random features.

<sup>5</sup>Choromanski, et al. (2016). Recycling randomness with structure for sublinear time kernel expansions.

<sup>6</sup>Avron, et al. (2017). Random Fourier features for kernel ridge regression: A approximation bounds and statistical guarantees.

Integral representation of kernel function

$$k(x, x') = \int_{\Omega} \psi(w, x) \psi(w, x') p(w) dw = \int_{\Omega} f_{xx'}(w) p(w) dw,$$

where  $p(w)$  is  $\mathcal{N}(0, \sigma_p^2 \mathbf{I})$  – density associated with the kernel and  $\psi(\cdot, x)$  is a feature map.

# Quadrature rules

1. Change variables to spherical-radial coordinates  
( $w = r\mathbf{z}$ ,  $\mathbf{z}^\top \mathbf{z} = 1$ ,  $w^\top w = r$ )

$$k(x, x') = \frac{(2\pi)^{-\frac{d}{2}}}{2} \int_{U_d} \int_{-\infty}^{\infty} e^{-\frac{r^2}{2}} |r|^{d-1} f_{xx'}(r\mathbf{z}) dr d\mathbf{z},$$

where  $U_d$  is a  $d$ -dimensional unit sphere.

2. Use stochastic radial-spherical rules.



- **Spherical-radial rules<sup>6</sup>** of degree  $(n, p)$

$$SR_{\mathbf{Q}, \rho}^{(n, p)}(f) = \sum_{j=1}^p \tilde{w}_j \sum_{i=1}^n \frac{w_i}{2} (f(-\rho_i \mathbf{Q} \mathbf{z}_j) + f(\rho_i \mathbf{Q} \mathbf{z}_j))$$

- Let  $p = 2n + 1$ . Then the weights are chosen such that the rule is exact for polynomials of degree  $2n + 1$ .

---

<sup>6</sup>Genz, A., & Monahan, J. (1998). Stochastic integration rules for infinite regions.

## Examples

► Degree 1 rule

$$SR_{\mathbf{Q},\rho}^{(1,1)}(f) = \frac{f_{xx'}(\rho\mathbf{Q}\mathbf{z}) + f_{xx'}(-\rho\mathbf{Q}\mathbf{z})}{2},$$

where

- $\rho \sim \chi(d)$ ,  $\mathbf{Q}$  – random orthogonal matrix,
- $\mathbf{z}$  – point on unit sphere.

→ Classical Random Fourier Features

► Degree (1,3) rule

$$SR_{\mathbf{Q},\rho}(f)^{(1,3)} = \sum_{i=1}^d \frac{f_{xx'}(-\rho\mathbf{Q}\mathbf{e}_i) + f_{xx'}(\rho\mathbf{Q}\mathbf{e}_i)}{2d},$$

where  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^\top$  with 1 in the  $i$ -th position.

→ Orthogonal Random Features!

## Degree (3, 3) rule

### Degree (3, 3) rule

$$k_{QBF}(x, x') = \left(1 - \frac{1}{d+1} \sum_{j=1}^{d+1} \frac{d}{\rho_j^2}\right) f_{xx'}(\mathbf{0}) + \frac{d}{d+1} \sum_{j=1}^{d+1} \left[ \frac{f_{xx'}(-\rho_j \mathbf{Q} \mathbf{v}_j) + f_{xx'}(\rho_j \mathbf{Q} \mathbf{v}_j)}{2\rho_j^2} \right],$$

where

- $\rho_j \sim \chi(d+2)$
- $\mathbf{v}_j$  is the  $j$ 'th vertex of unit  $d$ -simplex  $\mathbf{V}$
- $\mathbf{Q}$  is a random  $d \times d$  orthogonal matrix.

Use structured  $\mathbf{Q}$  matrix to speed up feature generation

# Empirical results

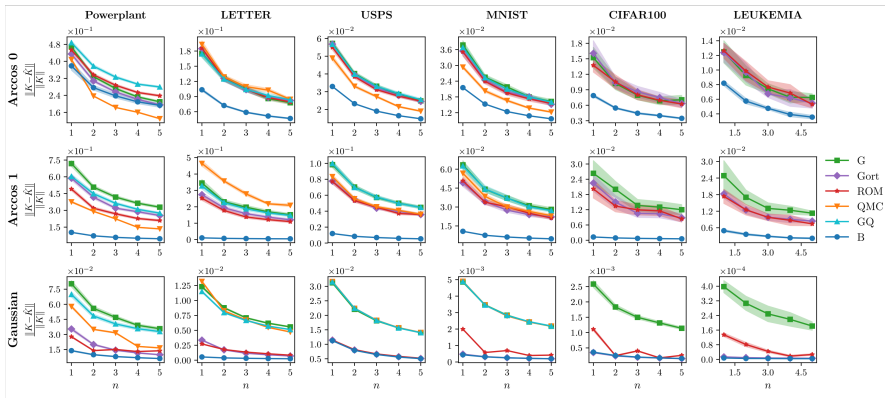


Figure: Error of kernel approximation on different datasets

# Empirical results

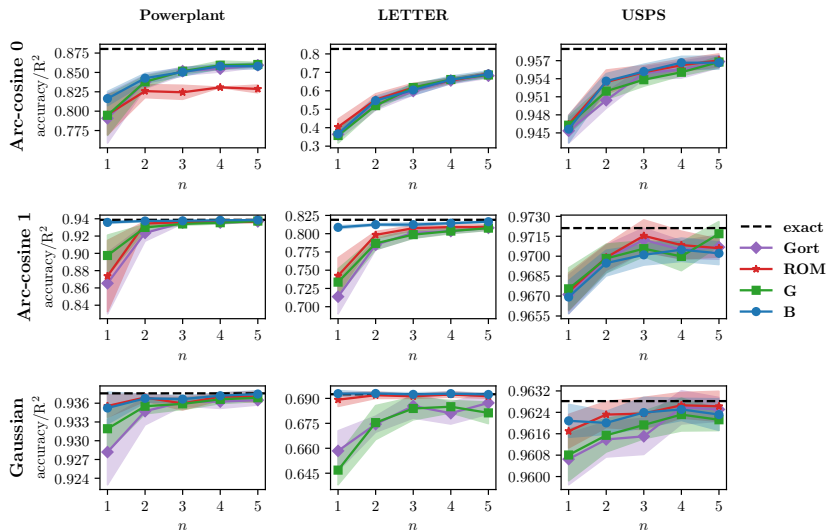


Figure: Error of regression/classification on different datasets

# Score Matching

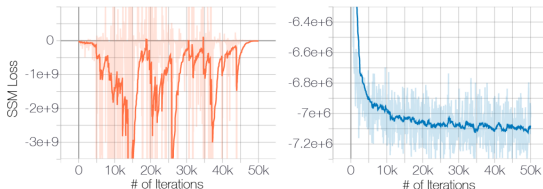
- ▶ Given a data set  $\{\mathbf{x}_i\}_{i=1}^n$ ,  $\mathbf{x}_i \sim p_0(\mathbf{x})$ , estimate unknown density  $p_0(\mathbf{x})$ .
- ▶ Find  $p_\theta(\mathbf{x}) \in \mathcal{P}$  that minimizes Fisher divergence

$$J(p_0 \| p_\theta) = \frac{1}{2} \int p_0(\mathbf{x}) \|\nabla \log p_\theta(\mathbf{x}) - \nabla \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}$$

- ▶ Equivalent **score matching** objective

$$J_{SM}(p_0 \| p_\theta) = \mathbb{E}_{p_0} \left[ \Delta \log p_\theta(\mathbf{x}) + \frac{1}{2} \|\nabla \log p_\theta(\mathbf{x})\|_2^2 \right]$$

- ▶ Need to compute second derivatives
- ▶  $\text{supp } p_0 \neq \text{supp } p_\theta$



Song, Y., Ermon, S. (2019). *Generative modeling by estimating gradients of the data distribution.*

**Figure:** Left: original and PCA reconstructions of images. Middle: score matching loss. Right: score matching loss with noisy data

# Kernel Exponential Family and Score Matching

- Consider distributions that satisfy

$$\log p_{\theta}(\mathbf{x}) = f(\mathbf{x}) + \log q_0(\mathbf{x}), \quad f \in \mathcal{H},$$

where  $\mathcal{H}$  is an RKHS with kernel  $k$ ,  $q_0$  is some generating density.

- Solution

$$f(\mathbf{x}) = -\frac{\xi}{\lambda} + \sum_{a=1}^n \sum_{i=1}^d \beta_{(\alpha-1)d+i} \partial_i k(\mathbf{x}_a, \cdot),$$

$$\xi = \frac{1}{n} \sum_{a=1}^n \sum_{i=1}^d \left( \partial_i k(\mathbf{x}_a, \cdot) \partial_i \log q_0(\mathbf{x}_a) + \partial_i^2 k(\mathbf{x}_a, \cdot) \right),$$

$$(\mathbf{G} + n\lambda\mathbf{I}) \beta = \frac{1}{\lambda} \mathbf{h},$$

$$\mathbf{G}_{(a-1)d+i, (b-1)d+j} = \partial_i \partial_{j+d} k(\mathbf{x}_a, \mathbf{x}_b),$$

$$\mathbf{h} = \langle \xi, \partial_i k(\mathbf{x}_a, \cdot) \rangle_{\mathcal{H}}$$



## Denoising score matching

Adding noise to input data is equivalent to convolution of the loss with noise distribution:

$$\begin{aligned} & \mathbb{E}_{p_\varepsilon} \mathbb{E}_{p_0} \left[ \Delta \log p_\theta(\mathbf{x} + \varepsilon) + \frac{1}{2} \|\nabla \log p_\theta(\mathbf{x} + \varepsilon)\|^2 \right] \\ &= \mathbb{E}_{p_0} \left[ \left( \left( \Delta \log p_\theta(\cdot) + \frac{1}{2} \|\nabla \log p_\theta(\cdot)\|^2 \right) * p_\varepsilon \right) (\mathbf{x}) \right] \end{aligned}$$

Solution derived only for Random Features

$$\hat{f}(\mathbf{x}) = \frac{1}{\lambda} \phi(\mathbf{x})^\top (\mathbf{H} + n\lambda \mathbf{I})^{-1} \mathbf{H} \mathbf{h} - \frac{1}{\lambda} \phi(\mathbf{x})^\top \mathbf{h},$$

where

$$\mathbf{H} = \int p_\varepsilon(\mathbf{y}) \partial \Phi_y^\top \partial \Phi_y d\mathbf{y}, \quad \mathbf{h} = \frac{1}{n} (\partial^2 \Phi_z * p(\mathbf{z}))^\top \mathbf{1},$$

and

$$[\partial \Phi_y]_{(a-1)d+i} = \partial_i \phi^\top(\mathbf{W}(\mathbf{x}_a + \mathbf{y}))$$

# Denoising Score Matching with Random Features

- ▶ Regularization:

$$\mathbf{h}_i \sim e^{-\sigma_\varepsilon^2 \|\mathbf{w}_i\|_2^2},$$

i.e. small weights for high-frequency components

- ▶ Explicit dependence on noise parameters  $\rightarrow$  easier to tune them
- ▶ We learn  $p_0 * p_\varepsilon \rightarrow$  trade-off between
  - ▶ stability of convergence
  - ▶ closeness to  $p_0$

# Experimental setup

- ▶ Adjust kernel parameters on training set using Denoising Score Matching
- ▶ Adjust noise variance on hold-out set using Score Matching
- ▶ Compare 3 models: proposed approach (**DSM RFF**), score matching with RFF (**SM RFF**), score matching with Nyström approximation (**Nyström**)
- ▶ Datasets
  - ▶ synthetic 2D data sets: Cosine, Uniform, Banana, Funnel, Rings
  - ▶ UCI: RedWine, WhiteWine, MiniBoone
- ▶ Metrics:
  - ▶ Log-likelihood
  - ▶ Wasserstein distance
  - ▶ Fisher divergence (for synthetic datasets)

# Mixtures

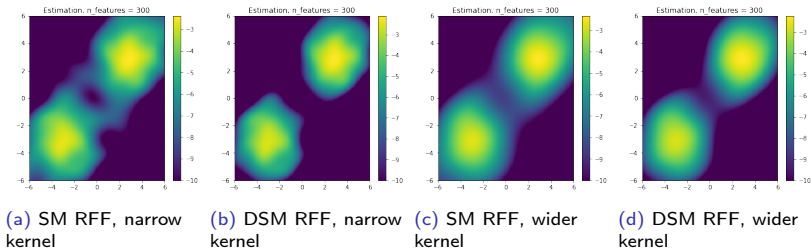


Figure: Mixture of Gaussians

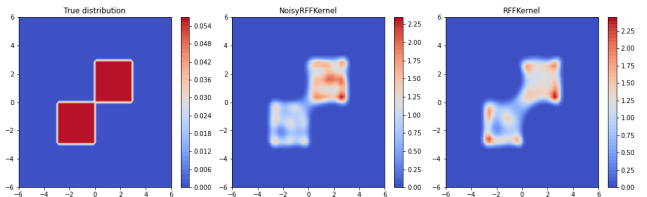


Figure: Mixture of Uniforms

# Experiments

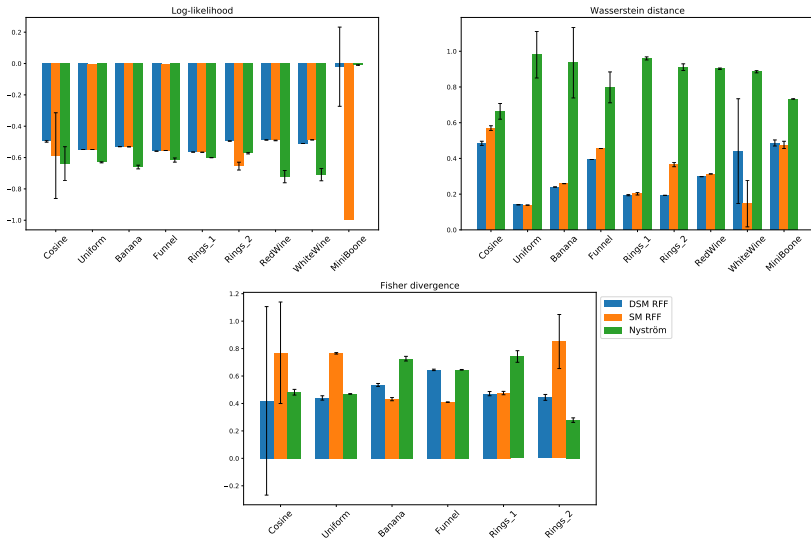


Figure: Comparison on synthetic and UCI data sets

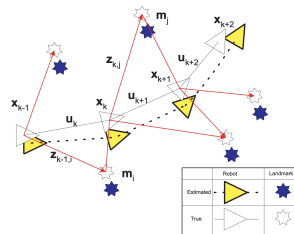
# Simultaneous Localization and Mapping (SLAM)

- ▶ Consider mobile robot moving in an environment.
- ▶ At each time step  $t_i$  we obtain measurements  $\mathbf{z}_i$

corresponding to some **landmarks**  $\mathbf{l} = \begin{bmatrix} \mathbf{l}_1 \\ \dots \\ \mathbf{l}_M \end{bmatrix}$

$$\mathbf{z}_i = \mathbf{h}(\mathbf{x}(t_i), \mathbf{l}) + \mathbf{n}_i, \quad \mathbf{n}_i \sim \mathcal{N}(0, \mathbf{R}_i).$$

- ▶ Control variables  $\mathbf{u}_i$  (maybe missing or given at different timestamps)
- ▶ We want to estimate both the robot trajectory  $\mathbf{x}(t_1), \dots, \mathbf{x}(t_T)$  and landmarks  $\mathbf{l}$ .
- ▶ Time-continuous SLAM: estimate trajectory as a function of time  $\mathbf{x}(t)$ .



*Durrant-Whyte, H., Bailey, T. (2006)*

► Assumptions

$$\begin{aligned} \mathbf{x}(t) &\sim \mathcal{GP}(\boldsymbol{\mu}_x, \mathbf{k}(t, t')) \\ \mathbf{l} &\sim \mathcal{N}(\boldsymbol{\mu}_l, \mathbf{L}) \end{aligned}$$

- Let us denote  $\boldsymbol{\theta} = \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{l} \end{bmatrix}$ . We want to maximize the posterior

$$p(\boldsymbol{\theta}|\mathbf{z}) \propto p(\mathbf{z}|\boldsymbol{\theta})p(\boldsymbol{\theta}) = -\frac{1}{2} \left( \sum_{i=1}^T \|\mathbf{z}_i - \mathbf{h}(\boldsymbol{\theta}(t_i))\|_{\mathbf{R}_i}^2 + \|\boldsymbol{\theta} - \boldsymbol{\mu}\|_{\mathbf{P}}^2 \right) = -J$$

- Random Features for locations

$$\mathbf{x}(t) = \begin{bmatrix} x(t) \\ y(t) \\ \alpha(t) \end{bmatrix} = \boldsymbol{\psi}(t)^\top \mathbf{b} + \varepsilon, \quad \begin{array}{l} \mathbf{b} \in \mathbb{R}^{D \times 3}, i = \overline{1, 3} \\ \mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_b, \mathbf{B}) \end{array}$$

where  $\boldsymbol{\psi}(t)$  – random features,  $\mathbf{B}$  is block-diagonal covariance matrix.

- Gauss-Newton method

$$(\delta \mathbf{b}^*, \delta \mathbf{l}) = \underset{\delta \mathbf{b}, \delta \mathbf{l}}{\operatorname{argmin}} \sum_{i=1}^T \left\| \mathbf{z}_i - \mathbf{h}(\boldsymbol{\psi}(t_i) \bar{\mathbf{b}}, \bar{\mathbf{l}}) - \mathbf{H}_i \Psi_i \delta \mathbf{b} \right\|_{\mathbf{R}_i}^2 + \left\| \bar{\mathbf{b}} + \delta \mathbf{b} - \boldsymbol{\mu}_b \right\|_{\mathbf{B}}^2 + \left\| \bar{\mathbf{l}} + \delta \mathbf{l} - \boldsymbol{\mu}_l \right\|_{\mathbf{L}}^2,$$

where  $\bar{\mathbf{b}}$  is current solution.



- ▶ A lot of works that use splines
- ▶ State-space model <sup>7</sup>

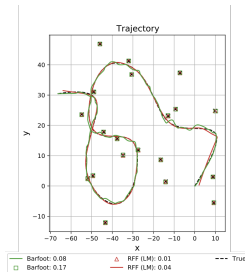
$$\dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{v}(t) + \mathbf{F}(t)\mathbf{w}(t),$$

- $\mathbf{A}(t), \mathbf{F}(t)$  are time-dependent system matrices,
  - $\mathbf{w}(t) \sim \mathcal{GP}(0, \mathbf{Q}_C \delta(t - t'))$
- ▶ Solution is GP with block-tridiagonal inverse  $\mathbf{K}$  matrix.
- ▶ Assumes Markovian-trajectories.
- ▶ RBF kernel can better in some cases (like noisy observations)

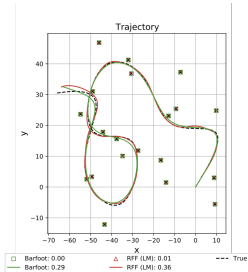
---

<sup>7</sup>Barfoot et al. (2014). Batch Continuous-Time Trajectory Estimation as Exactly Gaussian Process Regression

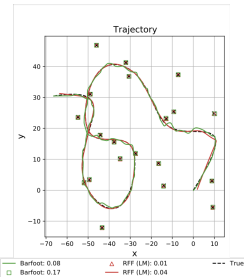
# Synthetic trajectories



(a) Bearing only



(b) Range only



(c) Range-Bearing

## Synthetic trajectories

		Pos.	Rot.	Landmarks
RangeBearing	RFF	<b>0.033</b>	0.0018	<b>0.019</b>
	Barfoot	0.096	0.0013	0.191
Range	RFF	0.309	0.0197	0.004
	Barfoot	<b>0.208</b>	0.0114	0.001
Bearing	RFF	<b>0.036</b>	0.0016	<b>0.018</b>
	Barfoot	0.096	0.0013	0.191

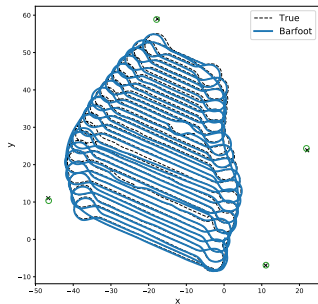
Table: Relative Pose Errors

Relative position errors don't take into account drift:

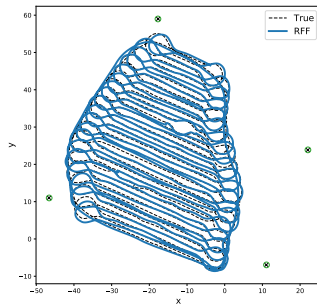
$$\text{RPE} = \sum_{i=1}^T \|\delta \hat{\mathbf{x}}_i \ominus \delta \mathbf{x}_i\|$$

# Autonomous Lawn-Mower

Range only data set



(a) State-Space Model



(b) Random Features based

# Summary

## Random Features

- ▶ proposed Quadrature-based Features
- ▶ accurate kernel approximation
- ▶ in downstream tasks benefit is smaller

## Score Matching

- ▶ exact solution for Denoising Score Matching with RFF was proposed
- ▶ Natural regularization
- ▶ Faster than Nyström-type approximation

## SLAM

- ▶ Random Features give dense covariance  $\rightarrow$  better accuracy in case of noisy data
- ▶ Random Features can oscillate