

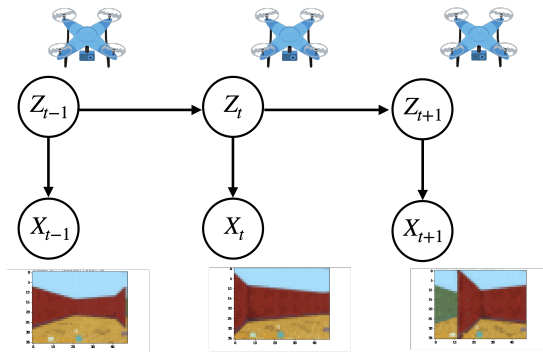
# Importance Weighting Approach in Kernel Bayes' Rule

Liyuan Xu <sup>1</sup>   Yutian Chen <sup>2</sup>  
Arnaud Doucet <sup>2</sup>   Arthur Gretton <sup>1</sup>

<sup>1</sup>Gatsby Unit

<sup>2</sup>DeepMind

## Drone Localization



- Predict drone location  $Z_t$  from camera images  $X_1; \dots; X_t$ .
- One approach: [Bayes' Filter](#)

## Bayes' Filter

$$P(Z_{t+1}jX_{1;\dots;t}) = \int P(Z_{t+1}jZ_t)dP(Z_tjX_{1;\dots;t})$$

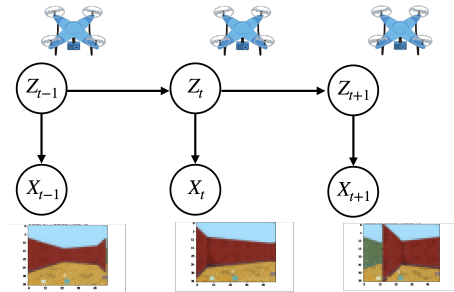
- Known as “sum rule” .

## Bayes' Filter

$$P(Z_{t+1} | X_{1, \dots, t+1}) / P(X_{t+1} | Z_{t+1}) P(Z_{t+1} | X_{1, \dots, t+1})$$

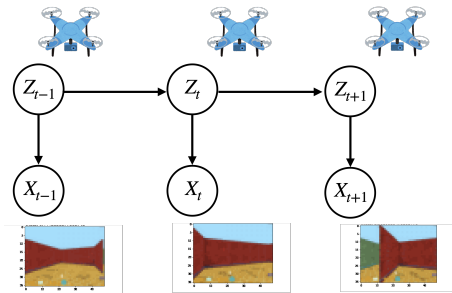
- Bayesian update of prior  $P(Z_{t+1} | X_{1, \dots, t})$  given observation  $X_{t+1}$
- Likelihood function is  $P(X_{t+1} | Z_{t+1}) = P(X_t | Z_t)$

## Difficulty of Bayes' Filter



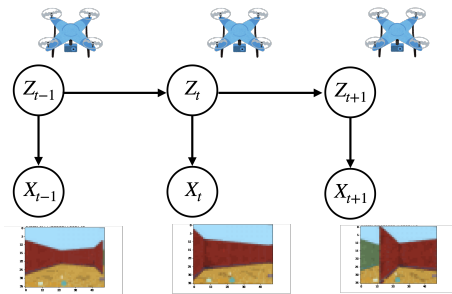
- There is **no explicit model** of  $P(Z_{t+1}|Z_t)$  or  $P(X_t|Z_t)$ .
  - ! Need to learn them from data  $fX_t; Z_tg$
- $P(Z_t|X_1, \dots, t)$  is assumed to be in a specific parametric form.
  - ! Might cause a bias in estimation.
- Desirable to use **non-parametric representation of distributions**.

## Difficulty of Bayes' Filter



- There is **no explicit model** of  $P(Z_{t+1}|Z_t)$  or  $P(X_t|Z_t)$ .  
! Need to learn them from data  $fX_t; Z_tg$
- $P(Z_t|X_{1;...;t})$  is assumed to be in a specific parametric form.  
! Might cause a bias in estimation.
- Desirable to use **non-parametric representation of distributions**.

## Difficulty of Bayes' Filter



- There is **no explicit model** of  $P(Z_{t+1}|Z_t)$  or  $P(X_t|Z_t)$ .  
! Need to learn them from data  $fX_t; Z_tg$
- $P(Z_t|X_{1;\dots;t})$  is assumed to be in a specific parametric form.  
! Might cause a bias in estimation.
- Desirable to use **non-parametric representation of distributions**.

## RKHS Embeddings

- Define kernel  $k(x; x')$  and accompanied feature map  $\phi(x)$ .

$$k(x; x') = \langle \phi(x); \phi(x') \rangle; \quad f(x) = hf; \langle \phi(x); \phi(x') \rangle$$

- Mean embedding  $\mu_P$  of distribution  $P$  is defined as

$$\mu_P = E_P[\phi(X)];$$

- It can be generalized to conditional distribution  $P_{X|Z}$

$$\mu_{P_{X|Z}}(z) = E_{P_{X|Z}}[\phi(X)|Z = z]$$



## RKHS Embeddings

- Define kernel  $k(x; x')$  and accompanied feature map  $\phi(x)$ .

$$k(x; x') = \langle \phi(x); \phi(x') \rangle; \quad f(x) = hf; \langle \phi(x); \phi(x') \rangle$$

- Mean embedding  $\mu_P$  of distribution  $P$  is defined as

$$\mu_P = E_P [ \phi(X) ];$$

- It can be generalized to conditional distribution  $P_{X|Z}$

$$\mu_{P_{X|Z}}(z) = E_{P_{X|Z}} [ \phi(X) | Z = z ]$$

## RKHS Embeddings

- Define kernel  $k(x; x')$  and accompanied feature map  $\phi(x)$ .

$$k(x; x') = \langle \phi(x); \phi(x') \rangle; \quad f(x) = hf; \langle \phi(x); \phi(x') \rangle$$

- Mean embedding  $\mu_P$  of distribution  $P$  is defined as

$$\mu_P = E_P [ \phi(X) ];$$

- It can be generalized to conditional distribution  $P_{X|Z}$

$$\mu_{P_{X|Z}}(z) = E_{P_{X|Z}} [ \phi(X) | Z = z ]$$

## Advantage of RKHS Embeddings

- Embedding  $\rho$  can uniquely determine the distribution.
- Embeddings can be non-parametrically estimated from  $fX_i; Z_i g_{i=1}^n$  as

$$P = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i)}; \quad P_{X|Z}(z) = \frac{1}{n} \sum_{i=1}^n w_i(z) \delta_{(X_i)}$$

for some weighting function  $w_i(z)$ .

- Kernel Bayes' Filter:  
Represent distributions  $P(Z_t | X_1, \dots, t)$  using embeddings

## Advantage of RKHS Embeddings

- Embedding  $P$  can uniquely determine the distribution.
- Embeddings can be non-parametrically estimated from  $\{X_i; Z_i\}_{i=1}^n$  as

$$P = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i)}; \quad P_{X|Z}(z) = \frac{1}{n} \sum_{i=1}^n w_i(z) \delta_{(X_i)}$$

for some weighting function  $w_i(z)$ .

- Kernel Bayes' Filter:

Represent distributions  $P(Z_t | X_1, \dots, X_t)$  using embeddings

## Advantage of RKHS Embeddings

- Embedding  $P$  can uniquely determine the distribution.
- Embeddings can be non-parametrically estimated from  $\{X_i; Z_i\}_{i=1}^n$  as

$$P = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i)}; \quad P_{X|Z}(z) = \frac{1}{n} \sum_{i=1}^n w_i(z) \delta_{(X_i)}$$

for some weighting function  $w_i(z)$ .

- Kernel Bayes' Filter:  
Represent distributions  $P(Z_t | X_{1:t})$  using embeddings

## Kernel Bayes' Filter

$$P_{Z_t | X_1, \dots, t} \quad \text{Data: } f_{Z_t; Z_{t+1}} \theta \quad P_{Z_{t+1} | X_1, \dots, t}$$

- Use kernel sum rule [Song et al. 2011] for  $P_{Z_{t+1} | X_1, \dots, X_t}$ .

## Kernel Bayes' Filter

$$P_{Z_{t+1}|X_1, \dots, t} \quad \text{Data: } f_{Z_t|X_t} g \quad \text{Observation: } X_{t+1} \quad P_{Z_{t+1}|X_1, \dots, t+1}$$

- Bayesian update given the **embedding** of the prior  $P_{Z_{t+1}|X_1, \dots, t}$
- This update is called **kernel Bayes' rule**.

## Kernel Bayes' Rule [Fukumizu+ 2013]

Given

- Training data  $\{X_i; Z_i\}_{i=1}^n$   $P(X|Z)P(Z)$
- Embedding of prior  $(Z)$

Outputs posterior embedding

$$Q(x) = \int_Z \frac{P(x|z) \phi(z)}{P(x|Z)} dz$$



## Contribution

Proposed a novel instance of kernel Bayes' rule.

- Based on **importance weighting**.
- Achieves **superior numerical stability** to existing work [Fukumizu+ 2013].
- Admits the use of **neural network feature** in kernel Bayes' rule.

## DeepMind Lab Experiment

- A drone is rotating in a maze.
- Latent  $Z_t$ : True angle of the drone.
- Observation  $X_t$ : The image observed at the noisy version of  $Z_t$
- Task: Predict  $Z_t$  from  $X_1; \dots; X_t$

## DeepMind Lab Experiment

- Original: Existing kernel Bayes rule [Fukumizu+ 2013]
- LSTM: Directly regress  $Z_t$  from  $X_{t-5:t}$
- IW: New kernel Bayes' rule with RKHS feature
- IW(NN): New kernel Bayes' rule with adaptive feature

## DeepMind Lab Experiment

- Original: Existing kernel Bayes rule [Fukumizu+ 2013]
- LSTM: Directly regress  $Z_t$  from  $X_{t-5:t}$
- IW: New kernel Bayes' rule with RKHS feature
- IW(NN): New kernel Bayes' rule with adaptive feature

## DeepMind Lab Experiment

- Original: Existing kernel Bayes rule [Fukumizu+ 2013]
- LSTM: Directly regress  $Z_t$  from  $X_{t-5:t}$
- IW: New kernel Bayes' rule with RKHS feature
- IW(NN): New kernel Bayes' rule with adaptive feature

## DeepMind Lab Experiment

- Original: Existing kernel Bayes rule [Fukumizu+ 2013]
- LSTM: Directly regress  $Z_t$  from  $X_{t-5:t}$
- IW: New kernel Bayes' rule with RKHS feature
- IW(NN): New kernel Bayes' rule with adaptive feature