# Bayesian Interpretations of RKHS Embedding Methods



#### David Duvenaud

Cambridge University Computational and Biological Learning Lab

December 8, 2012

# Outline

- Optimally-weighted Herding is Bayesian Quadrature
  - Kernel Herding
  - Bayesian Quadrature
  - Unifying Results
  - Demos
- Frequentist Methods, Bayesian Takeaways
  - Kernel Herding
  - Mean Embeddings
  - Kernel Two-sample Test
  - Hilbert-Schmidt Independence Criterion
  - Determinantal Point Processes

#### The Quadrature Problem

• We want to estimate an integral

$$Z = \int f(x)p(x)dx$$

- Most computational problems in Bayesian inference correspond to integrals:
  - Expectations
  - Marginal distributions
  - Integrating out nuisance parameters
  - Normalization constants



$$\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$



 Monte Carlo methods: Sample from p(x), take empirical mean:

$$\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$

Possibly sub-optimal for two reasons:



$$\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$

- Possibly sub-optimal for two reasons:
  - Random bunching up



$$\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$

- Possibly sub-optimal for two reasons:
  - Random bunching up
  - Often, nearby function values will be similar



$$\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$$

- Possibly sub-optimal for two reasons:
  - Random bunching up
  - Often, nearby function values will be similar
- Quasi-Monte Carlo methods spread out samples to achieve faster convergence.



# Kernel Herding [Welling et. al., 2009, Chen et. al., 2010]

• A sequential procedure for choosing sample locations, depending on previous locations.

Kernel Herding [Welling et. al., 2009, Chen et. al., 2010]

- A sequential procedure for choosing sample locations, depending on previous locations.
- Keeps estimate rule  $\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$

# Kernel Herding [Welling et. al., 2009, Chen et. al., 2010]

- A sequential procedure for choosing sample locations, depending on previous locations.
- Keeps estimate rule  $\hat{Z} = \frac{1}{N} \sum_{i=1}^{N} f(x_i)$
- Almost  $\mathcal{O}(1/N)$  convergence instead of  $\mathcal{O}(1/\sqrt{N})$  typical of random sampling, by spreading out samples.



#### Kernel Herding Objective

KH was found to minimize Maximum Mean Discrepancy:

$$\mathrm{MMD}_{\mathcal{H}}(p,q) = \sup_{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}} = 1}} \left| \int f(x)p(x)dx - \int f(x)q(x)dx \right|$$

#### Kernel Herding Objective

KH was found to minimize Maximum Mean Discrepancy:

$$\mathrm{MMD}_{\mathcal{H}}(p,q) = \sup_{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}} = 1}} \left| \int f(x)p(x)dx - \int f(x)q(x)dx \right|$$

In KH, p(x) is true distribution, and q(x) is a set of point masses at sample locations  $\{x_1, \ldots, x_N\}$ :

$$\epsilon_{KH}(\{x_1,\ldots,x_N\}) = \text{MMD}_{\mathcal{H}}\left(p,\underbrace{\frac{1}{N}\sum_{n=1}^N \delta_{x_n}}_{q(x)}\right)$$

# Kernel Herding

 Assuming function is in a Reproducing Kernel Hilbert Space defined by k(·, ·), MMD has closed form.

# Kernel Herding

- Assuming function is in a Reproducing Kernel Hilbert Space defined by k(·, ·), MMD has closed form.
- When sequentially minimizing MMD, new point is added at:

$$x_{N+1} = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$

$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



$$x_{N+1} = \arg_{x \in \mathcal{X}} \left[ 2 \int k(x, x') p(x') dx' - \frac{1}{N+1} \sum_{m=1}^{N} k(x, x_m) \right]$$



# Kernel Herding Summary

- A sequential sampling method which minimizes a worst-case divergence, given that f(x) belongs to a given RKHS.
- Like Monte Carlo, weights all samples  $f(x_s)$  equally when estimating Z:

$$\hat{Z} = \sum_{i=1}^{N} \frac{1}{N} f(x_i)$$
## Kernel Herding Summary

- A sequential sampling method which minimizes a worst-case divergence, given that f(x) belongs to a given RKHS.
- Like Monte Carlo, weights all samples  $f(x_s)$  equally when estimating Z:

$$\hat{Z} = \sum_{i=1}^{N} \frac{1}{N} f(x_i)$$

• What if we allowed different weights?

## Kernel Herding Summary

- A sequential sampling method which minimizes a worst-case divergence, given that f(x) belongs to a given RKHS.
- Like Monte Carlo, weights all samples  $f(x_s)$  equally when estimating Z:

$$\hat{Z} = \sum_{i=1}^{N} \frac{1}{N} f(x_i)$$

- What if we allowed different weights?
- [Bach et. al. 2012] looked at weighted herding strategies, showed improvement in convergence rates.

## Kernel Herding Summary

- A sequential sampling method which minimizes a worst-case divergence, given that f(x) belongs to a given RKHS.
- Like Monte Carlo, weights all samples  $f(x_s)$  equally when estimating Z:

$$\hat{Z} = \sum_{i=1}^{N} \frac{1}{N} f(x_i)$$

- What if we allowed different weights?
- [Bach et. al. 2012] looked at weighted herding strategies, showed improvement in convergence rates.

# Can we reason about the optimal weighting strategy?

- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



[O'Hagan 1987, Diaconis 1988, Rasmussen & Ghahramani 2003]

- Places a GP prior on f, defined by  $k(\cdot, \cdot)$  and a mean function.
- Posterior over *f* implies posterior over *Z*.



• Can choose samples however we want.

## Bayesian Quadrature Estimator

Posterior over Z has mean linear in  $f(x_s)$ :

$$\mathbb{E}_{\rm GP}\left[Z|f(x_s)\right] = \sum_{i=1}^N w_{BQ}^{(i)}f(x_i)$$

where

$$w_{BQ} = z^T K^{-1}$$
 and  $z_n = \int k(x, x_n) p(x) dx$ 

#### Bayesian Quadrature Estimator

Posterior over Z has mean linear in  $f(x_s)$ :

$$\mathbb{E}_{\rm GP}\left[Z|f(x_s)\right] = \sum_{i=1}^N w_{BQ}^{(i)}f(x_i)$$

where



• Natural to minimize the posterior variance of Z:

$$\mathbb{V}\left[Z|f(x_s)\right] = \iint k(x, x')p(x)p(x')dxdx' - z^T K^{-1}z$$
  
where  $z_n = \int k(x, x_n)p(x)dx$ 

• Natural to minimize the posterior variance of Z:

$$\mathbb{V}\left[Z|f(x_s)\right] = \iint k(x, x')p(x)p(x')dxdx' - z^T K^{-1}z$$
  
where  $z_n = \int k(x, x_n)p(x)dx$ 

• Favours samples in regions where p(x) is high, but where covariance with other sample locations is low. Similar flavour to herding objective.

• Natural to minimize the posterior variance of Z:

$$\mathbb{V}\left[Z|f(x_s)\right] = \iint k(x, x')p(x)p(x')dxdx' - z^T K^{-1}z$$
  
where  $z_n = \int k(x, x_n)p(x)dx$ 

- Favours samples in regions where p(x) is high, but where covariance with other sample locations is low. Similar flavour to herding objective.
- Does not depend on function values

• Natural to minimize the posterior variance of Z:

$$\mathbb{V}\left[Z|f(x_s)\right] = \iint k(x, x')p(x)p(x')dxdx' - z^T K^{-1}z$$
  
where  $z_n = \int k(x, x_n)p(x)dx$ 

- Favours samples in regions where p(x) is high, but where covariance with other sample locations is low. Similar flavour to herding objective.
- Does not depend on function values
- Can choose samples sequentially: Sequential Bayesian Quadrature.

## **Relating Objectives**

KH and BQ have completely different motivations:

- KH minimizes a worst-case bound
- BQ minimizes a posterior variance

Is there any correspondence?

## **Relating Objectives**

KH and BQ have completely different motivations:

- KH minimizes a worst-case bound
- BQ minimizes a posterior variance

Is there any correspondence?

#### First Main Result

$$\mathbb{V}\left[Z|f(x_s)\right] = \mathrm{MMD}^2(p, q_{\mathrm{BQ}})$$

Where

$$q_{\rm BQ}(x) = \sum_{n=1}^{N} w_{\rm BQ}^{(n)} \delta_{x_n}(x)$$

## **Relating Objectives**

KH and BQ have completely different motivations:

- KH minimizes a worst-case bound
- BQ minimizes a posterior variance

Is there any correspondence?

#### First Main Result

$$\mathbb{V}\left[Z|f(x_s)\right] = \mathrm{MMD}^2(p, q_{\mathrm{BQ}})$$

Where

$$q_{\rm BQ}(x) = \sum_{n=1}^{N} w_{\rm BQ}^{(n)} \delta_{x_n}(x)$$

#### BQ is minimizing KH objective

• KH and BQ are minimizing the same objective, but BQ has freedom to choose weights.

- KH and BQ are minimizing the same objective, but BQ has freedom to choose weights.
- How does this affect performance?

- KH and BQ are minimizing the same objective, but BQ has freedom to choose weights.
- How does this affect performance?

#### Second Main Result

BQ estimator is the optimal weighting strategy:

$$\mathbb{V}\left[Z|f(x_s)\right] = \inf_{w \in \mathbb{R}^N} \sup_{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}} \mathcal{H} = 1}} \left| \int f(x)p(x)dx - \sum_{n=1}^N w_n f(x_n) \right|^2$$

- KH and BQ are minimizing the same objective, but BQ has freedom to choose weights.
- How does this affect performance?

#### Second Main Result

BQ estimator is the optimal weighting strategy:

$$\mathbb{V}\left[Z|f(x_{s})\right] = \inf_{w \in \mathbb{R}^{N}} \sup_{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}} \mathcal{H} = 1}} \left| \int f(x)p(x)dx - \sum_{n=1}^{N} w_{n}f(x_{n}) \right|^{2}$$

 $\mathbb{V}[Z|f(x_s)]$  has two interpretations:

- Bayesian: posterior variance of Z under a GP prior.
- Frequentist: tight bound on estimation error of Z.

What is rate of convergence of BQ?

Expected Variance / MMD



What is rate of convergence of BQ?

Expected Variance / MMD



What is rate of convergence of BQ?

Expected Variance / MMD



What is rate of convergence of BQ?

Expected Variance / MMD

**Empirical Rates in RKHS** 





What is rate of convergence of BQ?

Expected Variance / MMD

## Empirical Rates out of RKHS



What is rate of convergence of BQ?

Expected Variance / MMD

Bound on Bayesian Error





• Posterior variance of Z under GP prior is equivalent to Maximum Mean Discrepancy.

- Posterior variance of Z under GP prior is equivalent to Maximum Mean Discrepancy.
- RKHS assumption gives a tight, closed-form upper bound on Bayesian error.

- Posterior variance of Z under GP prior is equivalent to Maximum Mean Discrepancy.
- RKHS assumption gives a tight, closed-form upper bound on Bayesian error.
- BQ has very fast, but unknown convergence rate.

- Posterior variance of Z under GP prior is equivalent to Maximum Mean Discrepancy.
- RKHS assumption gives a tight, closed-form upper bound on Bayesian error.
- BQ has very fast, but unknown convergence rate.
- The optimal weighted herding strategy is Bayesian quadrature.

- Posterior variance of Z under GP prior is equivalent to Maximum Mean Discrepancy.
- RKHS assumption gives a tight, closed-form upper bound on Bayesian error.
- BQ has very fast, but unknown convergence rate.
- The optimal weighted herding strategy is Bayesian quadrature.
- Joint work with Ferenc Huzsar
# Outline

- Optimally-weighted Herding is Bayesian Quadrature
  - Kernel Herding
  - Bayesian Quadrature
  - Unifying Results
  - Demos
- Frequentist Methods, Bayesian Takeaways
  - Kernel Herding
  - Mean Embeddings
  - Kernel Two-sample Test
  - Hilbert-Schmidt Independence Criterion
  - Determinantal Point Processes











Takeaway: Herding assumptions are innapropriate for inference



Mean Embedding Interpretation

• [Muandet & Ghahramani, 2012] showed that if  $f \sim {
m GP}$ ,

## Mean Embedding Interpretation

• [Muandet & Ghahramani, 2012] showed that if  $f \sim {
m GP}$ ,

$$\mu_{p(x)} = \int \phi(x)p(x)dx$$
  
=  $\int k(x, \cdot)p(x)dx$   
=  $\mathbb{E}_{f\sim GP} \left[ \int f(x)f(\cdot)p(x)dx \right]$   
=  $\mathbb{E}_{f\sim GP} \left[ f(\cdot) \int f(x)p(x)dx \right]$   
=  $\operatorname{cov}_{f\sim GP} (f(\cdot), Z_p)$ 

#### Mean Embedding Interpretation

• [Muandet & Ghahramani, 2012] showed that if  $f \sim {
m GP}$ ,

$$\mu_{p(x)} = \int \phi(x)p(x)dx$$
  
=  $\int k(x, \cdot)p(x)dx$   
=  $\mathbb{E}_{f\sim GP} \left[ \int f(x)f(\cdot)p(x)dx \right]$   
=  $\mathbb{E}_{f\sim GP} \left[ f(\cdot) \int f(x)p(x)dx \right]$   
=  $\operatorname{cov}_{f\sim GP} (f(\cdot), Z_p)$ 

 μ<sub>p(x)</sub> is the covariance the function with its integral with respect to p(x).

• How to test whether two distributions p(x) and q(x) are the same?

- How to test whether two distributions p(x) and q(x) are the same?
- New test statistic: MMD(p, q) [Gretton et. al, 2005]

- How to test whether two distributions p(x) and q(x) are the same?
- New test statistic: MMD(p, q) [Gretton et. al, 2005]
- Equivalent Bayesian intepretation:

$$\mathrm{MMD}_{k}^{2}(p,q) = \mathbb{V}_{f \sim \mathrm{GP}_{k}}\left[\int f(x)p(x)dx - \int f(x)q(x)dx\right]$$

p and q are similar if integrals of functions drawn from a GP prior have similar integrals.

- How to test whether two distributions p(x) and q(x) are the same?
- New test statistic: MMD(p, q) [Gretton et. al, 2005]
- Equivalent Bayesian intepretation:

$$\mathrm{MMD}_{k}^{2}(p,q) = \mathbb{V}_{f \sim \mathrm{GP}_{k}}\left[\int f(x)p(x)dx - \int f(x)q(x)dx\right]$$

p and q are similar if integrals of functions drawn from a GP prior have similar integrals.

Possible takeaways: Decision-theoretic choice of kernel, sampling-based methods for computing MMD

• Given samples  $\{X, Y\} \sim p(x, y)$ , how to test whether p(x, y) = p(x)p(y)?

- Given samples  $\{X, Y\} \sim p(x, y)$ , how to test whether p(x, y) = p(x)p(y)?
- New test statistic based on infinite-dimensional Frobenius norm of cross-covariance matrix of features of x and y: [Gretton et. al, 2005]

- Given samples  $\{X, Y\} \sim p(x, y)$ , how to test whether p(x, y) = p(x)p(y)?
- New test statistic based on infinite-dimensional Frobenius norm of cross-covariance matrix of features of x and y: [Gretton et. al, 2005]

$$\begin{aligned} &\text{HSIC}(p(x, y), k_x, k_y) = ||C_{xy}||_{HS}^2 \\ &= \mathbb{E}_{x, x', y, y'} \left[ k_x(x, x') k_y(y, y) \right] + \mathbb{E}_{x, x'} \left[ k_x(x, x') \right] \mathbb{E}_{y, y'} \left[ k_y(y, y') \right] \\ &- 2 \mathbb{E}_{x, y} \left[ \mathbb{E}_{x'} \left[ k_x(x, x') \right] \mathbb{E}_{y'} \left[ k_y(y, y') \right] \right] \end{aligned}$$

- Given samples  $\{X, Y\} \sim p(x, y)$ , how to test whether p(x, y) = p(x)p(y)?
- New result: Assuming  $k(x, y, x', y') = k_x(x, x')k_y(y, y')$

$$HSIC(p(x, y), k_x, k_y) =$$
  
=  $\mathbb{V}_{f \sim GP_k} \left[ \int f(x, y) p(x, y) dx dy - \int f(x, y) p(x) p(y) dx dy \right]$ 

• Probability of a set  $P(\mathcal{X}) = |K(\mathcal{X}, \mathcal{X})|$ 

- Probability of a set  $P(\mathcal{X}) = |\mathcal{K}(\mathcal{X}, \mathcal{X})|$
- Greedy MAP maximizes  $P(\mathcal{X} \cup x_i)$

- Probability of a set  $P(\mathcal{X}) = |\mathcal{K}(\mathcal{X}, \mathcal{X})|$
- Greedy MAP maximizes  $P(\mathcal{X} \cup x_i)$

$$P(\mathcal{X} \cup x_i)$$

$$= |K(\mathcal{X} \cup x_i, \mathcal{X} \cup x_i|)$$

$$= |K(\mathcal{X}, \mathcal{X})| [k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)]$$

$$\propto k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)$$

$$= \mathbb{V}_{f \sim GP_k} [f(x_i)|\mathcal{X}]$$

- Probability of a set  $P(\mathcal{X}) = |K(\mathcal{X}, \mathcal{X})|$
- Greedy MAP maximizes  $P(\mathcal{X} \cup x_i)$

$$P(\mathcal{X} \cup x_i)$$

$$= |K(\mathcal{X} \cup x_i, \mathcal{X} \cup x_i|)$$

$$= |K(\mathcal{X}, \mathcal{X})| [k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)]$$

$$\propto k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)$$

$$= \mathbb{V}_{f \sim GP_k} [f(x_i)|\mathcal{X}]$$

 New DPP Point added at location with highest marginal variance in GP posterior, conditioned on the other points

- Probability of a set  $P(\mathcal{X}) = |K(\mathcal{X}, \mathcal{X})|$
- Greedy MAP maximizes  $P(\mathcal{X} \cup x_i)$

$$P(\mathcal{X} \cup x_i)$$

$$= |K(\mathcal{X} \cup x_i, \mathcal{X} \cup x_i|)$$

$$= |K(\mathcal{X}, \mathcal{X})| [k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)]$$

$$\propto k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)$$

$$= \mathbb{V}_{f \sim GP_k} [f(x_i)|\mathcal{X}]$$

- New DPP Point added at location with highest marginal variance in GP posterior, conditioned on the other points
- Related to sensor placement work by Andreas Krause

- Probability of a set  $P(\mathcal{X}) = |K(\mathcal{X}, \mathcal{X})|$
- Greedy MAP maximizes  $P(\mathcal{X} \cup x_i)$

$$P(\mathcal{X} \cup x_i)$$

$$= |K(\mathcal{X} \cup x_i, \mathcal{X} \cup x_i|)$$

$$= |K(\mathcal{X}, \mathcal{X})| [k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)]$$

$$\propto k(x_i, x_i) - k(x_i, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1}k(\mathcal{X}, x_i)$$

$$= \mathbb{V}_{f \sim GP_k} [f(x_i)|\mathcal{X}]$$

- New DPP Point added at location with highest marginal variance in GP posterior, conditioned on the other points
- Related to sensor placement work by Andreas Krause
- Thanks to Le Song and Roman Garnett

\_

Frequentist Method	Bayesian Method
Kernel herding	Bayesian Quadrature
Mean Embedding	Covariance with integral
Kernel Two-sample test	Variance of difference of integrals
HSIC	Variance of difference of integrals
DPP MAP	Marginal uncertainty in GP posterior

Frequentist Method	Bayesian Method
Kernel herding	Bayesian Quadrature
Mean Embedding	Covariance with integral
Kernel Two-sample test	Variance of difference of integrals
HSIC	Variance of difference of integrals
DPP MAP	Marginal uncertainty in GP posterior

Some possible extensions:

Frequentist Method	Bayesian Method
Kernel herding	Bayesian Quadrature
Mean Embedding	Covariance with integral
Kernel Two-sample test	Variance of difference of integrals
HSIC	Variance of difference of integrals
DPP MAP	Marginal uncertainty in GP posterior

Some possible extensions:

• Log-kernel herding for inference

Frequentist Method	Bayesian Method
Kernel herding	Bayesian Quadrature
Mean Embedding	Covariance with integral
Kernel Two-sample test	Variance of difference of integrals
HSIC	Variance of difference of integrals
DPP MAP	Marginal uncertainty in GP posterior

#### Some possible extensions:

- Log-kernel herding for inference
- Interpretation of conditional mean embeddings

Bayesian Method
Bayesian Quadrature
Covariance with integral
Variance of difference of integrals
Variance of difference of integrals
Marginal uncertainty in GP posterior

Some possible extensions:

- Log-kernel herding for inference
- Interpretation of conditional mean embeddings
- Different low-rank approximations based on sparse GPs

Frequentist Method	Bayesian Method
Kernel Regression	GP Regression
Functional ANOVA,	Additive Gaussian Processes
Hierchical Kernel Learning	
Kernel Bayes' Rule	Bayesian Quadrature for Ratios
RKHS Embeddings of	GP Dynamics Models
Conditional Distributions	
Kernel Message Passing	???

Frequentist Method	Bayesian Method
Kernel Regression	GP Regression
Functional ANOVA,	Additive Gaussian Processes
Hierchical Kernel Learning	
Kernel Bayes' Rule	Bayesian Quadrature for Ratios
RKHS Embeddings of	GP Dynamics Models
Conditional Distributions	
Kernel Message Passing	???



Thanks!