

Métodos Numéricos Probabilísticos

Roberto Díaz Morales

DTSC, Universidad Carlos III de Madrid



Índice

1. Introducción

2. Cuadratura

3. Muestreo

4. Cuadratura Bayesiana

5. Aprendizaje activo de la evidencia usando cuadratura bayesiana.

6. Conclusiones

Introducción

- Una de las tareas del aprendizaje máquina consiste en la inferencia sobre datos complejos.
- Dicha tarea incluye el desarrollo y análisis de métodos numéricos basados en teoría probabilística.
- Muchas de estas tareas se pueden ver como problemas de aprendizaje.
- Dos de las tareas más importantes son:
 - Cuadratura
 - Muestreo

Cuadratura

- ▶ Muchos modelos complejos requieren integrales computacionalmente intratables, con lo que han de ser aproximadas.
- ▶ En particular, muchos problemas de inferencia requieren integrar sobre funciones de probabilidad.

$$Z = \langle \ell \rangle = \int \ell(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

- $p(\mathbf{x})$ podría ser un “posterior” y $f(\mathbf{x})$ las etiquetas de nuestras muestras.
 - $p(\mathbf{x})$ podría ser un “prior” y $f(\mathbf{x})$ una verosimilitud.
 - ...
- ▶ Es complicado cuando se trabaja con grandes conjuntos de datos donde evaluar la verosimilitud sobre todo el conjunto de datos es muy costoso computacionalmente.

Índice

1. Introducción

2. Cuadratura

3. Muestreo

4. Cuadratura Bayesiana

5. Aprendizaje activo de la evidencia usando cuadratura bayesiana.

6. Conclusiones

Estimador de Monte Carlo

- ▶ Si $p(\mathbf{x})$ es una función de densidad de probabilidad y podemos obtener muestras de ella tenemos la aproximación de Monte Carlo:

$$\Phi = \langle \phi(\mathbf{x}) \rangle \equiv \int P(\mathbf{x}) \phi(\mathbf{x}) d^N \mathbf{x}$$

$$\hat{\Phi} \equiv \frac{1}{R} \sum_r \phi(\mathbf{x}^{(r)}) \quad \sigma^2 = \int d^N \mathbf{x} P(\mathbf{x}) (\phi(\mathbf{x}) - \Phi)^2$$

- ▶ Su varianza decae en un orden $O(1/R)$
- ▶ La principal objeción es que tener muestras que representen bien $P(\mathbf{x})$ no garantiza que representen bien $\Phi(\mathbf{x})$

Índice

1. Introducción

2. Cuadratura

3. Muestreo

4. Cuadratura Bayesiana

5. Aprendizaje activo de la evidencia usando cuadratura bayesiana.

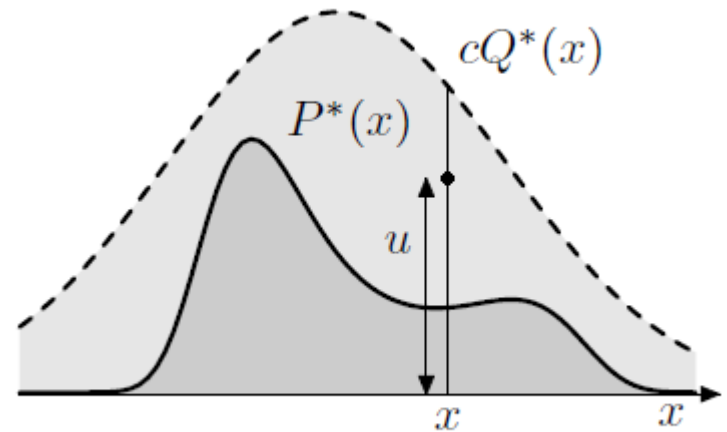
6. Conclusiones

Métodos de muestreo

- ▶ Existen multitud de métodos de muestreo cuando podemos evaluar $p(x)$:

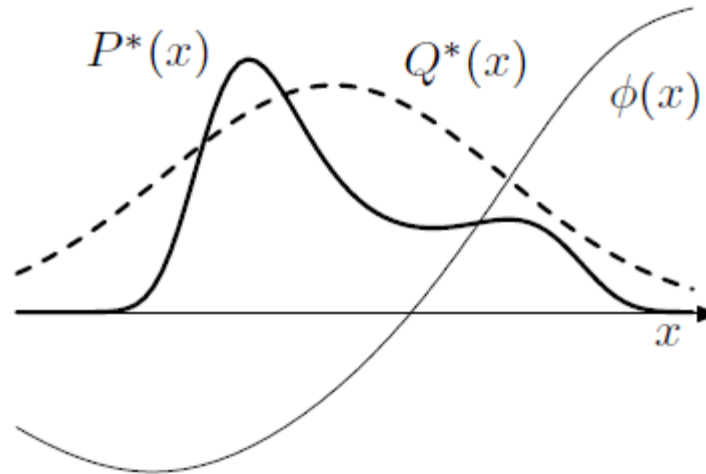
- Rejection Sampling:

- Se conoce $cQ(x) > P(x)$
- Se genera x_i muestra de $cQ(x)$
- Se genera $u \sim U(0, cQ(x_i))$
- Si $u < P(x)$ guardamos x_i
- Si $u > P(x)$ descartamos x_i
- Eficiencia según el parecido de $cQ(x)$ con $P(x)$



Métodos de muestreo

- Importance Sampling:



$$w_r \equiv \frac{P^*(x^{(r)})}{Q^*(x^{(r)})}$$

$$\hat{\Phi} \equiv \frac{\sum_r w_r \phi(x^{(r)})}{\sum_r w_r}$$

Métodos de muestreo MCMC

▶ Markov Chain Monte Carlo:

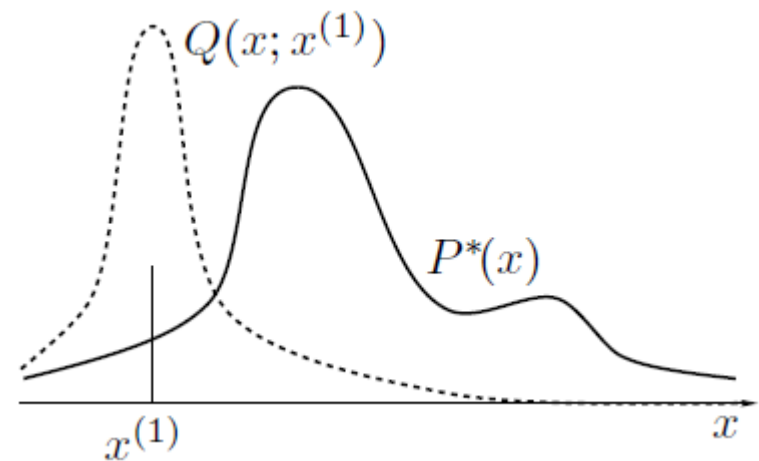
- Son técnicas de muestreo en las que la siguiente muestra depende de la muestra actual (estado) y de unas probabilidades de transición.

- Metropolis–Hasting:

- Se utiliza una función $Q(x)$ que depende de la muestra actual.
- Se evalúa a:

$$a = \frac{P^*(x')}{P^*(x^{(t)})} \frac{Q(x^{(t)}; x')}{Q(x'; x^{(t)})}$$

- Si $a > 1$ la muestra se acepta.
- Si $a < 1$ se acepta con probabilidad a .
- Si se acepta es el nuevo estado.
- Se debe dejar correr $T \simeq (L/\epsilon)^2$



Métodos de muestreo MCMC

- Gibbs–Sampling:
 - La siguiente muestra se obtiene con la distribución conjunta para cada una de las dimensiones.

$$x_1^{(t+1)} \sim P(x_1 | x_2^{(t)}, x_3^{(t)}, \dots, x_K^{(t)})$$

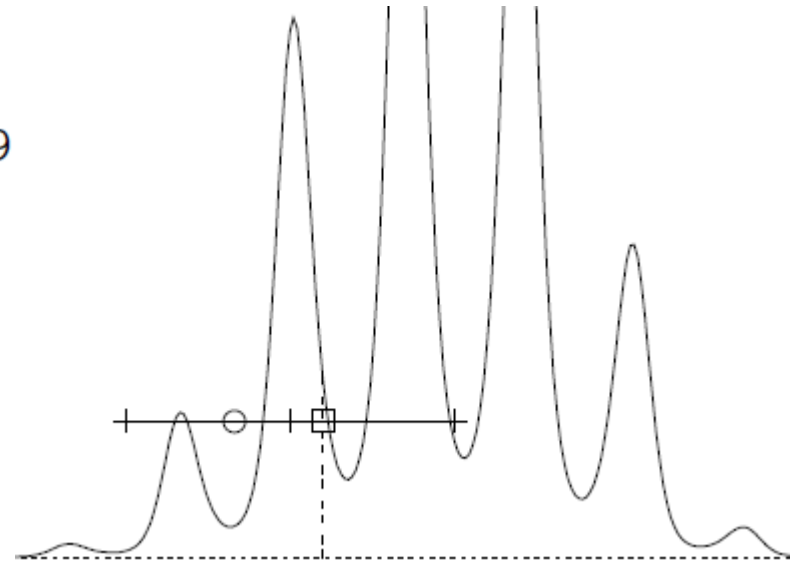
$$x_2^{(t+1)} \sim P(x_2 | x_1^{(t+1)}, x_3^{(t)}, \dots, x_K^{(t)})$$

$$x_3^{(t+1)} \sim P(x_3 | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_K^{(t)})$$

Métodos de muestreo MCMC

- Slice-Sampling:

- 1: evaluate $P^*(x)$
- 2: draw a vertical coordinate $u' \sim \text{Uniform}(0, P^*(x))$
- 3: create a horizontal interval (x_l, x_r) enclosing x
- 4: loop {
- 5: draw $x' \sim \text{Uniform}(x_l, x_r)$
- 6: evaluate $P^*(x')$
- 7: if $P^*(x') > u'$ break out of loop 4-9
- 8: else modify the interval (x_l, x_r)
- 9: }



Índice

1. Introducción

2. Cuadratura

3. Muestreo

4. Cuadratura Bayesiana

5. Aprendizaje activo de la evidencia usando cuadratura bayesiana.

6. Conclusiones

Cuadratura Bayesiana

- Rasmussen, C. E., & Ghahramani, Z.
Bayesian monte carlo. *Advances in neural information processing systems*, 15, 489-496.
- Dado un conjunto de muestras \mathcal{D} y realizando inferencia sobre f con \mathcal{D} , la media sobre funciones es la esperanza de $f(x)$ media.

$$\mathcal{D} = \{(x^{(i)}, f(x^{(i)})) | i = 1 \dots n\}$$

$$E_{f|\mathcal{D}}[\bar{f}_p] = \iint f(x)p(x)dx p(f|\mathcal{D})df$$

$$V_{f|\mathcal{D}}[\bar{f}_p] = \int \left[\int f(x)p(x)dx - \int \bar{f}(x')p(x')dx' \right]^2 p(f|\mathcal{D})df$$

Cuadratura Bayesiana

- Se agrupan los términos para un GP:

$$\begin{aligned} E_{f|\mathcal{D}}[\bar{f}_p] &= \iint f(x)p(x)dx p(f|\mathcal{D})df \\ &= \int \left[\int f(x)p(f|\mathcal{D})df \right] p(x)dx = \int \bar{f}_{\mathcal{D}}(x)p(x)dx \end{aligned}$$

$$\begin{aligned} V_{f|\mathcal{D}}[\bar{f}_p] &= \int \left[\int f(x)p(x)dx - \int \bar{f}(x')p(x')dx' \right]^2 p(f|\mathcal{D})df \\ &= \iiint [f(x) - \bar{f}(x)] [f(x') - \bar{f}(x')] p(f|\mathcal{D})df p(x)p(x')dx dx' \\ &= \iint \text{Cov}_{\mathcal{D}}(f(x), f(x')) p(x)p(x')dx dx', \end{aligned}$$

$$\bar{f}_{\mathcal{D}}(x) = k(x, \mathbf{x})K^{-1}\mathbf{f}, \quad \text{and} \quad \text{Cov}_{\mathcal{D}}(f(x), f(x')) = k(x, x') - k(x, \mathbf{x})K^{-1}k(\mathbf{x}, x')$$

Cuadratura Bayesiana

- En el caso general, introducir la formulación de los GP en la integral lleva a expresiones que son difíciles de evaluar, pero hay casos especiales:
 - Si $p(x)$ y la función de covarianza son ambas gaussianas, se obtienen expresiones analíticas (utilizando la cuadratura Bayes-Hermite):

$$E_{f|\mathcal{D}}[\bar{f}_p] = z^\top K^{-1} \mathbf{f}$$

$$V_{f|\mathcal{D}}[\bar{f}_p] = w_0 |2A^{-1}B + I|^{-1/2} - z^\top K^{-1} z$$

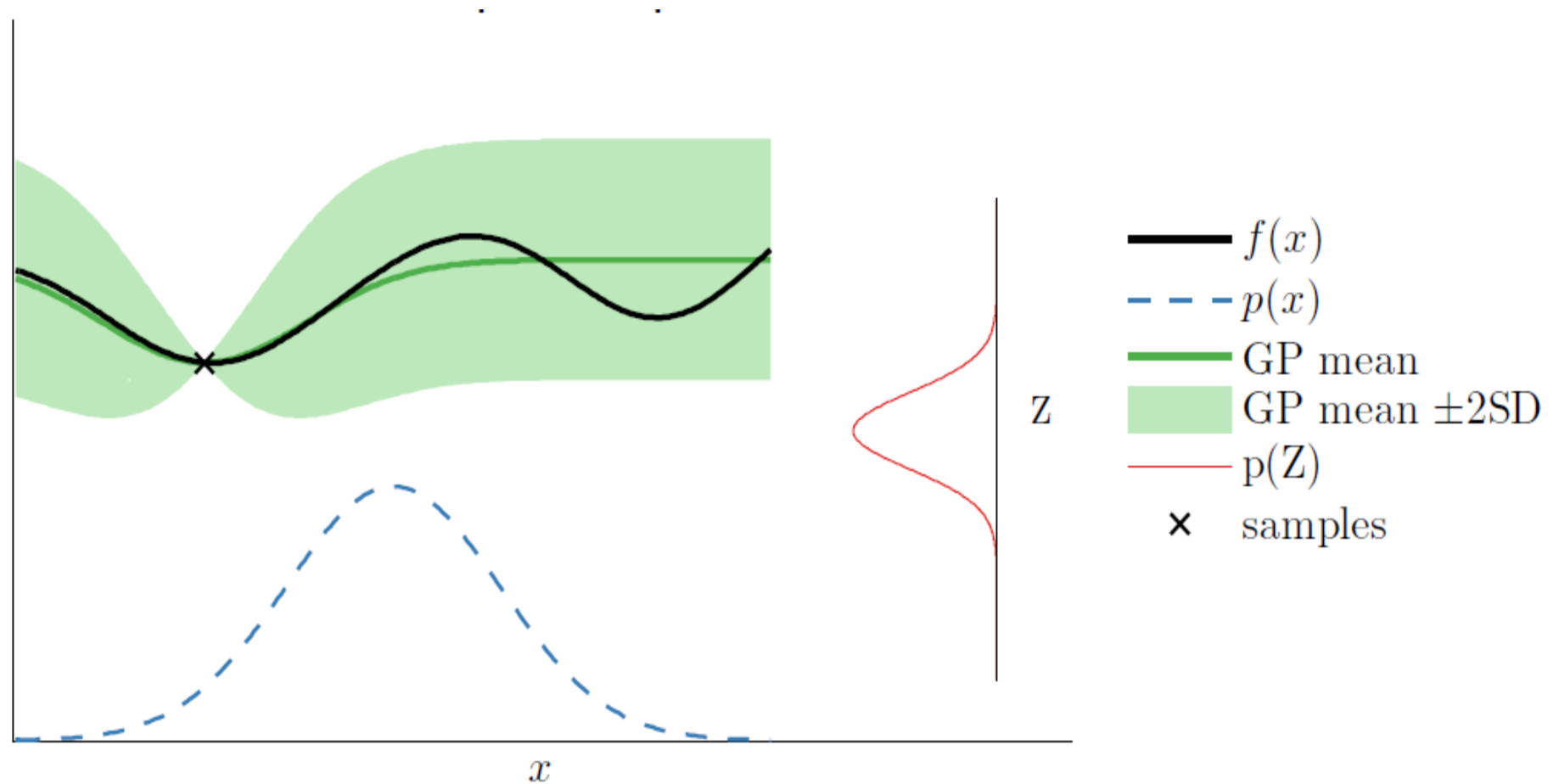
$$z = w_0 |A^{-1}B + I|^{-1/2} \exp[-0.5(a-b)^\top (A+B)^{-1} (a-b)]$$

$$p(x) = \mathcal{N}(b, B)$$

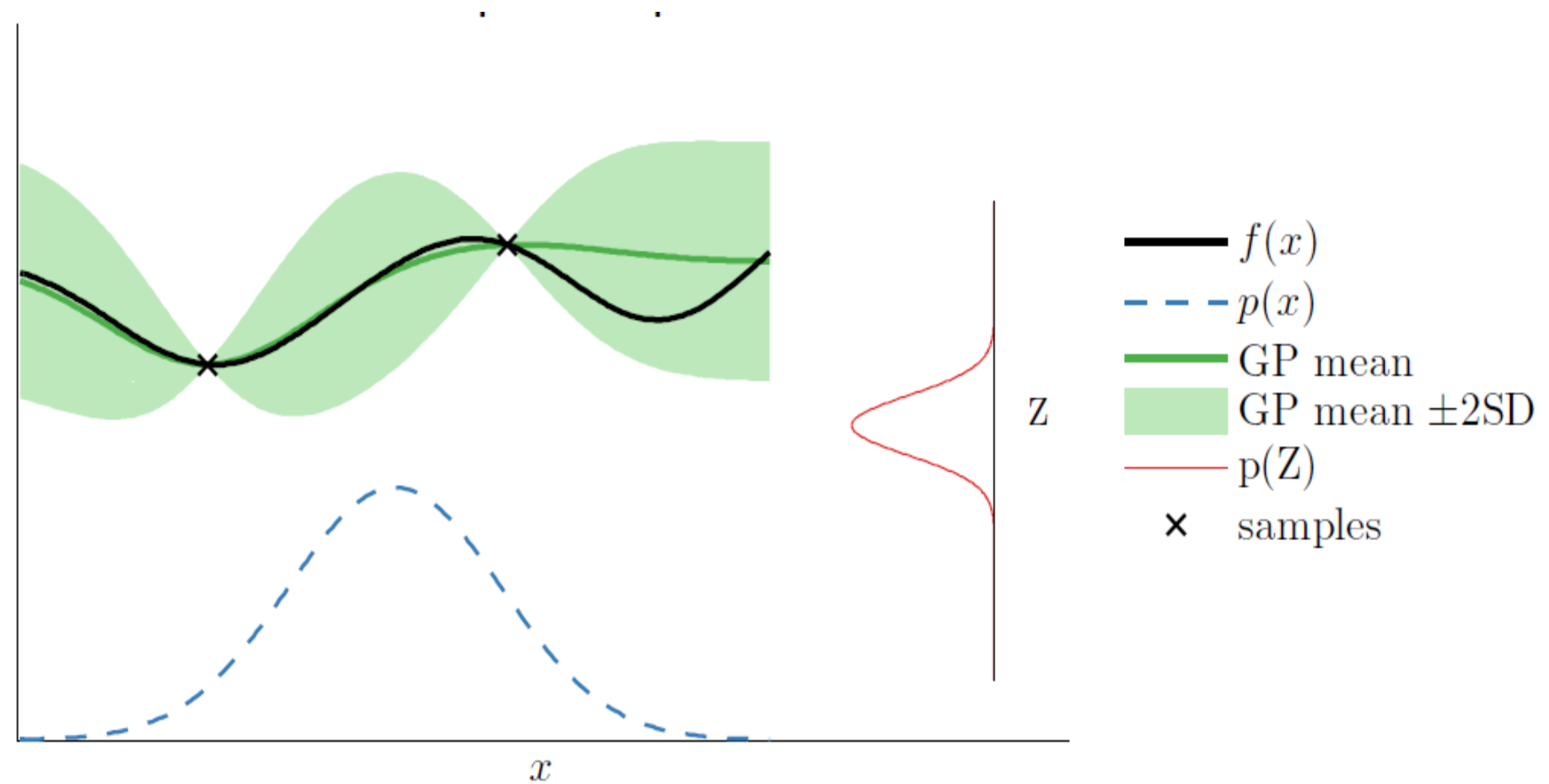
$$a_i = x^{(i)}$$

$$K_{pq} = \text{Cov}(f(x^{(p)}), f(x^{(q)})) = w_0 \exp\left(-\frac{1}{2} \sum_{d=1}^D (x_d^{(p)} - x_d^{(q)})^2 / w_d^2\right)$$

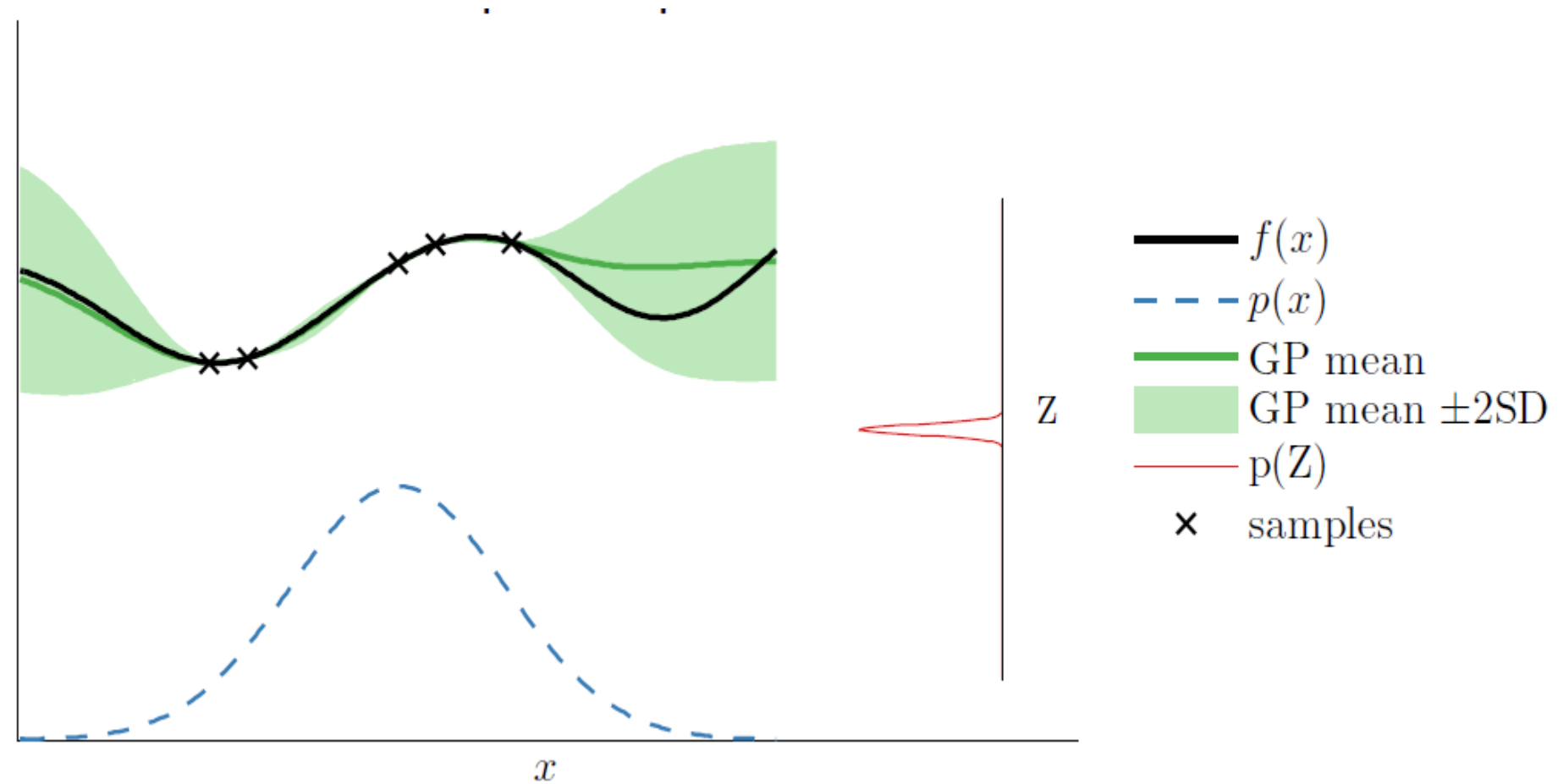
Cuadratura Bayesiana



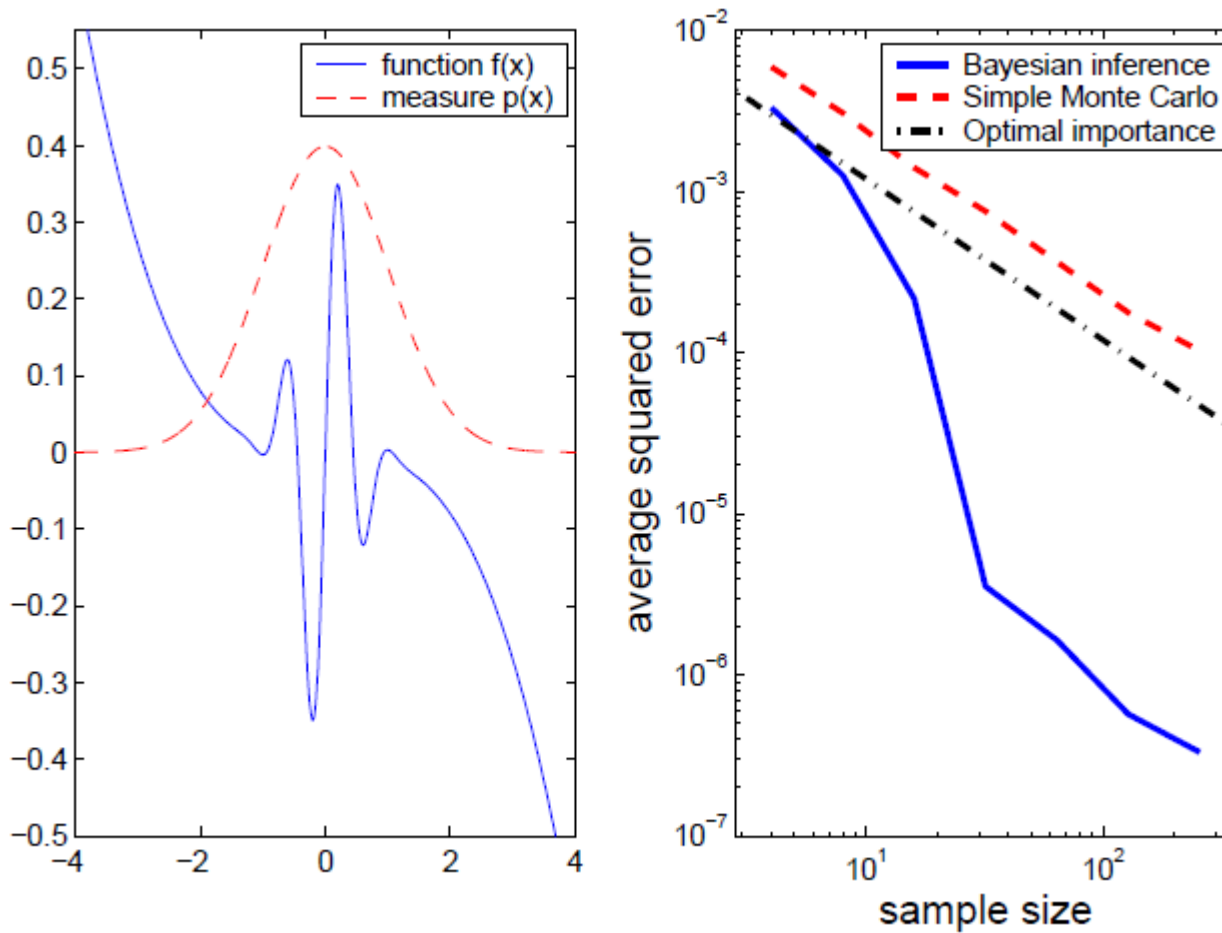
Cuadratura Bayesiana



Cuadratura Bayesiana



Cuadratura Bayesiana



Índice

1. Introducción

2. Cuadratura

3. Muestreo

4. Cuadratura Bayesiana

5. Aprendizaje activo de la evidencia usando cuadratura bayesiana.

6. Conclusiones

Aprendizaje Activo de Evidencia

- Osborne, M. A., Duvenaud, D., Garnett, R., Rasmussen, C. E., Roberts, S. J., & Ghahramani, Z. Active Learning of Model Evidence Using Bayesian Quadrature. NIPS 2012
- Se desea realizar una integral sobre una verosimilitud no negativa:

$$Z = \langle \ell \rangle = \int \ell(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

$$m(Z | \log \ell_s) = \int \left(\int \exp(\log \ell(x)) p(x) dx \right) \mathcal{N}(\log \ell; m_{\log \ell | s}, C_{\log \ell | s}) d \log \ell$$

Aprendizaje Activo de Evidencia

- Para poder tratar la integral, “lineariza” el problema.

$$\exp(\log \ell(x)) \simeq \exp(\log \ell_0(x)) + \exp(\log \ell_0(x)) (\log \ell(x) - \log \ell_0(x))$$

- Quedando como problema a resolver:

$$m(Z|\log \ell_s) \simeq m(Z|\log \ell_0, \log \ell_s) := \int \ell_0(x)p(x) dx + \int \ell_0(x)\Delta_{\log \ell|s}(x)p(x) dx$$

$$\Delta_{\log \ell|s} := m_{\log \ell|s} - \log \ell_0$$

Aprendizaje Activo de la Evidencia

- Como L_0 utiliza un GP estándar.
 - **$L_0 = m L | s$**
 - Prior con media 0
 - Covarianza Gaussiana
 - Utiliza el conjunto de datos X_s
- Utiliza un GP diferente para modelar $\Delta \log L | s$
 - Utiliza prior con media 0
 - Covarianza Gaussiana
 - Para entrenar utiliza X_s y datos aleatorios en hiper-
elipses alrededor de los puntos X_s

Aprendizaje Activo de Evidencia

- Quedando finalmente la media compuesta por dos términos analíticos.

$$m(Z|\log \ell_s) \simeq m(Z|\log \ell_0, \log \ell_s, \Delta_c) = m(Z|\ell_s) + m(\langle \ell \Delta_{\log \ell|s} \rangle | \ell_s, \Delta_c)$$

- Y una varianza que puede emplearse como diagnóstico de convergencia:

$$V(Z|\log \ell_0, \log \ell_s, \Delta_c) = S(Z|\log \ell_0, \log \ell_s) - m(Z|\log \ell_0, \log \ell_s, \Delta_c)^2$$

$$S(Z|\log \ell_0, \log \ell_s) := m(\langle \ell C_{\log \ell|s} \ell \rangle | \log \ell_s) + m(Z|\log \ell_0, \log \ell_s, \Delta_c)^2$$

$$\begin{aligned} V(Z|\log \ell_0, \log \ell_s, \Delta_c) &= m(\langle \ell C_{\log \ell|s} \ell \rangle | \log \ell_s) \\ &:= \iint m_{\ell|s}(x) m_{\ell|s}(x') C_{\log \ell|s}(x, x') p(x) p(x') dx dx' \end{aligned}$$

Aprendizaje Activo de Evidencia

- Aprendizaje activo:
 - Ya no es necesario coger muestras que pertenezcan a $p(x)$.
 - Cuando se han fijado los hyperparámetros, la varianza depende de la posición de las muestras escogidas.
 - Selecciona muestras que minimizan la varianza esperada:

$$x_a = \operatorname{argmin}_{x_a} \langle V(Z | \log \ell_0, \log \ell_{s,a}) \mid \log \ell_0, \log \ell_s \rangle$$

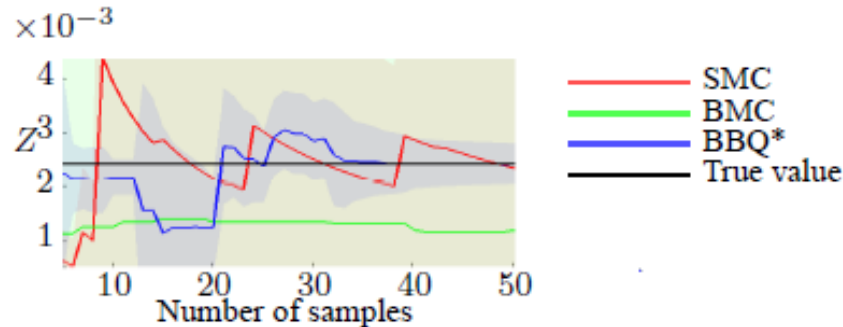
$$\begin{aligned} \langle V(Z | \log \ell_0, \log \ell_{s,a}) \mid \log \ell_0, \log \ell_s \rangle &= S(Z | \log \ell_0, \log \ell_s) - \int m(Z | \log \ell_0, \log \ell_{a,s}, \Delta_c)^2 \\ &\quad \times \mathcal{N}\left(\log \ell_a; \hat{m}_a, \hat{C}_a + \frac{\partial \hat{m}_a}{\partial w} C_w \frac{\partial \hat{m}_a^\top}{\partial w}\right) d\log \ell_a \end{aligned}$$

$$\hat{m}_a := m(\log \ell_a | \log \ell_s, \hat{w})$$

$$\hat{C}_a := V(\log \ell_a | \log \ell_s, \hat{w})$$

Aprendizaje Activo de Evidencia

- Resultados:



$$\text{ALE} := \frac{1}{N} \sum_{i=1}^N |\log m(Z_i) - \log Z_i|$$

Synthetic Results

Method	ALE
SMC	1.101
AIS	1.695
BMC	2.695
BO	6.760
BBQ	0.919

Real Results

Method	ALE
SMC	0.632
AIS	2.146
BMC	1.455
BO	0.635
BBQ	0.400

Referencias

- www.probablistic-numeric.org
- MacKay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge university press.
- Rasmussen, C. E., & Ghahramani, Z. Bayesian monte carlo. *Advances in neural information processing systems*, 15, 489-496.
- Osborne, M. A., Duvenaud, D., Garnett, R., Rasmussen, C. E., Roberts, S. J., & Ghahramani, Z. Active Learning of Model Evidence Using Bayesian Quadrature. NIPS 2012