

Efficient Mixture Modeling with RKHS Embeddings

A PAC-Bayesian Analysis

Matthew Higgs¹

¹Center for Computational Statistics and Machine Learning
University College London

Foundations and New Trends of PAC Bayesian Learning

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching
- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2^{nd} Order U -Statistics
 - Do we have a U -statistic?
- 3 Close
 - Choosing KL
 - Conclusion

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching
- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2nd Order U -Statistics
 - Do we have a U -statistic?
- 3 Close
 - Choosing KL
 - Conclusion

Maximum Mean Discrepancy

- $P, Q \in \mathcal{M}_+^1(\mathbb{R}^d)$, probability measures on \mathbb{R}^d .
- $\mathcal{F} \subset \mathbb{R}^{\mathbb{R}^d}$, measurable.
- $f_Q := \int_{\mathbb{R}^d} f(x) dQ(x)$.

Maximum Mean Discrepancy

- $P, Q \in \mathcal{M}_+^1(\mathbb{R}^d)$, probability measures on \mathbb{R}^d .
- $\mathcal{F} \subset \mathbb{R}^{\mathbb{R}^d}$, measurable.
- $f_Q := \int_{\mathbb{R}^d} f(x) dQ(x)$.

Definition (Maximum Mean Discrepancy [Gretton et al., 2008])

$$MMD_{\mathcal{F}}(Q, P) := \sup_{f \in \mathcal{F}} |f_Q - f_P|. \quad (1)$$

Maximum Mean Discrepancy

- $P, Q \in \mathcal{M}_+^1(\mathbb{R}^d)$, probability measures on \mathbb{R}^d .
- $\mathcal{F} \subset \mathbb{R}^{\mathbb{R}^d}$, measurable.
- $f_Q := \int_{\mathbb{R}^d} f(x) dQ(x)$.

Definition (Maximum Mean Discrepancy [Gretton et al., 2008])

$$MMD_{\mathcal{F}}(Q, P) := \sup_{f \in \mathcal{F}} |f_Q - f_P|. \quad (1)$$

Examples

- $\mathcal{F} := C_b(\mathbb{R}^d)$
- $\mathcal{F} := \{f : \|f\|_{\infty} \leq 1\}$
- $\mathcal{F} := \{\mathbf{1}_{(-\infty, t]}(\cdot) : t \in \mathbb{R}^d\}$,
- $\mathcal{F} := \{e^{i\langle \omega, \cdot \rangle} : \omega \in \mathbb{R}^d\}$,

Unit-ball in RKHS

- $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$, sym. p.d., $\sup_{x,x'} k(x, x') = C_k < \infty$.
- \mathcal{H} s.t. $\langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} = k(x, x')$, $\|\cdot\|_{\mathcal{H}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathcal{H}}}$.
- $k_Q(\cdot) := \int k(x, \cdot) dQ(x), \in \mathcal{H}$.
- $k_{Q,P} := \int \int k(x', x) dQ(x) dP(x'), \in \mathbb{R}$.

Unit-ball in RKHS

- $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$, sym. p.d., $\sup_{x,x'} k(x, x') = C_k < \infty$.
- \mathcal{H} s.t. $\langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} = k(x, x')$, $\|\cdot\|_{\mathcal{H}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathcal{H}}}$.
- $k_Q(\cdot) := \int k(x, \cdot) dQ(x), \in \mathcal{H}$.
- $k_{Q,P} := \int \int k(x', x) dQ(x) dP(x'), \in \mathbb{R}$.

Proposition ([Song et al., 2008])

For any $P, Q \in \mathcal{M}_+^1(\mathbb{R}^d)$, $\mathcal{F} = \{f : \|f\|_{\mathcal{H}} \leq 1\}$,

$$MMD_{\mathcal{F}}(Q, P) = \|k_Q - k_P\|_{\mathcal{H}}.$$

Unit-ball in RKHS

- $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$, sym. p.d., $\sup_{x,x'} k(x, x') = C_k < \infty$.
- \mathcal{H} s.t. $\langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} = k(x, x')$, $\|\cdot\|_{\mathcal{H}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathcal{H}}}$.
- $k_Q(\cdot) := \int k(x, \cdot) dQ(x), \in \mathcal{H}$.
- $k_{Q,P} := \int \int k(x', x) dQ(x) dP(x'), \in \mathbb{R}$.

Proposition ([Song et al., 2008])

For any $P, Q \in \mathcal{M}_+^1(\mathbb{R}^d)$, $\mathcal{F} = \{f : \|f\|_{\mathcal{H}} \leq 1\}$,

$$MMD_{\mathcal{F}}(Q, P) = \|k_Q - k_P\|_{\mathcal{H}}.$$

Corollary

$$MMD_{\mathcal{F}}^2(Q, P) = k_{Q,Q} - 2k_{Q,P} + k_{P,P}. \quad (2)$$

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching
- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2nd Order U -Statistics
 - Do we have a U -statistic?
- 3 Close
 - Choosing KL
 - Conclusion

Reproducing KMM Mixture Model

- $S := (x_1, \dots, x_n) \sim D^n$, $D_n := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.
- $\Theta_{mult} := \{\theta \in \mathbb{R}^l : \theta^\top \mathbf{1} = 1, \theta \succ \mathbf{0}\}$.
- Mixture $Q_\theta = \sum_{i=1}^l \theta_i Q_i$, $\theta \in \Theta_{mult}$, $Q_i(\mathbb{R}^d) = 1$.
- $\mathbf{R}(i, j) := k_{Q_i, Q_j}$, $\mathbf{L}(i) := k_{Q_i, D_n}$.

Reproducing KMM Mixture Model

- $S := (x_1, \dots, x_n) \sim D^n$, $D_n := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.
- $\Theta_{mult} := \{\theta \in \mathbb{R}^l : \theta^\top \mathbf{1} = 1, \theta \succ \mathbf{0}\}$.
- Mixture $Q_\theta = \sum_{i=1}^l \theta_i Q_i$, $\theta \in \Theta_{mult}$, $Q_i(\mathbb{R}^d) = 1$.
- $\mathbf{R}(i, j) := k_{Q_i, Q_j}$, $\mathbf{L}(i) := k_{Q_i, D_n}$.

Corollary (KMM [Song et al., 2008])

Minimisation of $\|k_{Q_\theta} - k_{D_n}\|_{\mathcal{H}}^2$ with reg. $\lambda \|\theta\|^2$, becomes QP

$$\theta_{\min} := \arg \min_{\theta \in \Theta_{mult}} \frac{1}{2} \theta^\top (\mathbf{R} + \lambda \mathbf{I}) \theta - \mathbf{L} \theta \quad (3)$$

Reproducing KMM Mixture Model

- $S := (x_1, \dots, x_n) \sim D^n$, $D_n := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.
- $\Theta_{mult} := \{\theta \in \mathbb{R}^l : \theta^\top \mathbf{1} = 1, \theta \succ \mathbf{0}\}$.
- Mixture $Q_\theta = \sum_{i=1}^l \theta_i Q_i$, $\theta \in \Theta_{mult}$, $Q_i(\mathbb{R}^d) = 1$.
- $\mathbf{R}(i, j) := k_{Q_i, Q_j}$, $\mathbf{L}(i) := k_{Q_i, D_n}$.

Corollary (KMM [Song et al., 2008])

Minimisation of $\|k_{Q_\theta} - k_{D_n}\|_{\mathcal{H}}^2$ with reg. $\lambda \|\theta\|^2$, becomes QP

$$\theta_{\min} := \arg \min_{\theta \in \Theta_{mult}} \frac{1}{2} \theta^\top (\mathbf{R} + \lambda \mathbf{I}) \theta - \mathbf{L} \theta \quad (3)$$

Question

What PAC-Bayesian statistical guarantees do we have for θ_{\min} ?

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching
- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2nd Order U -Statistics
 - Do we have a U -statistic?
- 3 Close
 - Choosing KL
 - Conclusion

U -Statistic PAC-Bayes Bound

- $S = (x_1, \dots, x_n) \sim D^n$.
- $(\{h\})$ measurable, $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $a \leq h(x, x') \leq b$.
- $\bar{U}_S(h) := \frac{1}{n(n-1)} \sum_{i \neq j} \frac{h(x_i, x_j) - a}{b - a}$ with mean $\bar{U}(h)$.
- $\text{kl}(q \| p) := q \ln \frac{q}{p} + (1 - q) \ln(1 - q) \frac{1 - q}{1 - p}$.
- $\Phi_C(q, p) := \ln \frac{1}{1 - [1 - e^{-C}]p} - C \cdot q$, $C \in \mathbb{R}$.

U -Statistic PAC-Bayes Bound

- $S = (x_1, \dots, x_n) \sim D^n$.
- $(\{h\})$ measurable, $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $a \leq h(x, x') \leq b$.
- $\bar{U}_S(h) := \frac{1}{n(n-1)} \sum_{i \neq j} \frac{h(x_i, x_j) - a}{b-a}$ with mean $\bar{U}(h)$.
- $\text{kl}(q||p) := q \ln \frac{q}{p} + (1-q) \ln(1-q) \frac{1-q}{1-p}$.
- $\Phi_C(q, p) := \ln \frac{1}{1-[1-e^{-C}]p} - C \cdot q$, $C \in \mathbb{R}$.

Corollary

For any $\pi \in \mathcal{M}_+^1(\{h\})$, with probability $\geq 1 - \delta$ over the draw of S , for any $\rho \in \mathcal{M}_+^1(\{h\})$

$$\text{kl}\left(\mathbb{E}_{h \sim \rho} \bar{U}_S(h) \middle| \middle| \mathbb{E}_{h \sim \rho} \bar{U}(h)\right) \leq \frac{1}{\lfloor n/2 \rfloor} \left[\text{KL}(\rho || \pi) + \ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta} \right], \quad (4)$$

$$\Phi_C\left(\mathbb{E}_{h \sim \rho} \bar{U}_S(h), \mathbb{E}_{h \sim \rho} \bar{U}(h)\right) \leq \frac{1}{\lfloor n/2 \rfloor} \left[\text{KL}(\rho || \pi) + \ln \frac{1}{\delta} \right]. \quad (5)$$

Proof: iid blocks

- $\bar{h}(\cdot) = \frac{h(\cdot) - a}{b - a}$.
- Permutations σ on $\{1, \dots, n\}$.
- $\bar{B}_{S, \sigma}(h) := \frac{1}{\lfloor n/2 \rfloor} \sum_{i=1}^{\lfloor n/2 \rfloor} \bar{h}(X_{\sigma(i)}, X_{\sigma(\lfloor n/2 \rfloor + i)})$.
- $\mathcal{D} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ convex on $[0, 1]^2$.

Proof: iid blocks

- $\bar{h}(\cdot) = \frac{h(\cdot) - a}{b - a}$.
- Permutations σ on $\{1, \dots, n\}$.
- $\bar{B}_{S, \sigma}(h) := \frac{1}{\lfloor n/2 \rfloor} \sum_{i=1}^{\lfloor n/2 \rfloor} \bar{h}(X_{\sigma(i)}, X_{\sigma(\lfloor n/2 \rfloor + i)})$.
- $\mathcal{D} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ convex on $[0, 1]^2$.

Theorem

For any $\pi \in \mathcal{M}_+^1(\{h\})$, with probability $\geq 1 - \delta$ over the draw of S , for any $\rho \in \mathcal{M}_+^1(\{h\})$

$$\mathcal{D}\left(\mathbb{E}_{h \sim \rho} \bar{U}_S(h), \mathbb{E}_{h \sim \rho} \bar{U}(h)\right) \leq \frac{1}{\lfloor n/2 \rfloor} \left[\text{KL}(\rho \| \pi) + \ln \frac{\mathcal{L}_{\mathcal{D}}}{\delta} \right], \quad (6)$$

where $\mathcal{L}_{\mathcal{D}} = \mathbb{E}_{h \sim \rho} \mathbb{E}_S e^{\lfloor n/2 \rfloor \left[\mathcal{D}\left(\mathbb{E}_{h \sim \rho} \bar{U}_S(h), \mathbb{E}_{h \sim \rho} \bar{U}(h)\right) \right]}$, and (by Jensen's)

$$\mathcal{L}_{\mathcal{D}} \leq \frac{1}{m!} \sum_{\sigma} \mathbb{E}_{h \sim \rho} \mathbb{E}_S e^{\lfloor n/2 \rfloor \left[\mathcal{D}\left(\bar{B}_{S, \sigma}(h), \mathbb{E}_S \bar{B}_{S, \sigma}(h)\right) \right]}. \quad (7)$$

Proof: Bounded to Bernoulli

[Maurer, 2004].

$X = (X_1, \dots, X_n)$ iid, $X_i \in [0, 1]$; $X' = (X'_1, \dots, X'_n)$ iid, $X'_i \in \{0, 1\}$,
 $\mathbb{E}X' = \mathbb{E}X$. Then $\phi : [0, 1]^n \rightarrow \mathbb{R}$ convex, then

$$\mathbb{E}\phi(X) \leq \mathbb{E}\phi(X'). \quad (8)$$



Proof: Bounded to Bernoulli

[Maurer, 2004].

$X = (X_1, \dots, X_n)$ iid, $X_i \in [0, 1]$; $X' = (X'_1, \dots, X'_n)$ iid, $X'_i \in \{0, 1\}$, $\mathbb{E}X' = \mathbb{E}X$. Then $\phi : [0, 1]^n \rightarrow \mathbb{R}$ convex, then

$$\mathbb{E}\phi(X) \leq \mathbb{E}\phi(X'). \quad (8)$$

□

Proposition

For $\mathcal{L}_{\mathcal{D}} \leq \frac{1}{m!} \sum_{\sigma} \mathbb{E}_{h \sim P} \mathbb{E}_{\mathcal{S}} e^{\lfloor n/2 \rfloor} [\mathcal{D}(\bar{B}_{S, \sigma}(h), \mathbb{E}_{\mathcal{S}} \bar{B}_{S, \sigma}(h))]$, $\mathcal{S} = (x_1, \dots, x_n)$,

$$\mathcal{L}_{\text{kl}} \leq \xi(\lfloor n/2 \rfloor) \quad (9)$$

$$\mathcal{L}_{\Phi_{\mathcal{C}}} \leq 1, \quad \forall \mathcal{C} \in \mathbb{R}^+ \quad (10)$$

where $\xi(m) = \mathcal{O}(\sqrt{m})$.

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching

- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2nd Order U -Statistics
 - Do we have a U -statistic?

- 3 Close
 - Choosing KL
 - Conclusion

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Answer

No. But almost...

- $V(\theta) := \frac{1}{n} \sum_{k=1}^n \|k_{Q_{\theta}} - k(x_k, \cdot)\|_{\mathcal{H}}^2$ (bias term).

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Answer

No. But almost...

- $V(\theta) := \frac{1}{n} \sum_{k=1}^n \|k_{Q_\theta} - k(x_k, \cdot)\|_{\mathcal{H}}^2$ (bias term).
- $\text{kl}_{1/2}(q||p) := \text{kl}(\frac{1}{2}q + \frac{1}{2}||\frac{1}{2}p + \frac{1}{2})$.
- NB: $Q_\theta = \sum_{i=1}^I \theta_i Q_i$.

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Answer

No. But almost...

- $V(\theta) := \frac{1}{n} \sum_{k=1}^n \|k_{Q_{\theta}} - k(x_k, \cdot)\|_{\mathcal{H}}^2$ (bias term).
- $\text{kl}_{1/2}(q||p) := \text{kl}\left(\frac{1}{2}q + \frac{1}{2}||\frac{1}{2}p + \frac{1}{2}\right)$.
- NB: $Q_{\theta} = \sum_{i=1}^I \theta_i Q_i$.

Theorem

For any RKHS \mathcal{H} with kernel k , $\mathcal{F} = \{f : \|f\|_{\mathcal{H}} \leq 1\}$, any $\pi \in \mathcal{M}_+^1(\Theta_{\text{mult}})$, with probability $\geq 1 - \delta$ over the draw of S , for any $\hat{\rho} \in \mathcal{M}_+^1(\Theta_{\text{mult}})$ with mean $\hat{\theta} \in \Theta_{\text{mult}}$:

$$\text{kl}_{1/2} \left(\frac{n \|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2 - V(\hat{\theta})}{2C_k(n-1)} \left\| \frac{\|k_{Q_{\hat{\theta}}} - k_D\|_{\mathcal{H}}^2}{2C_k} \right. \right) \leq \frac{2\text{KL}(\hat{\rho}||\pi) + \ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta}}{\lfloor n/2 \rfloor}.$$

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Proof.

- $h_{(\theta, \theta')}(x, x') := \langle k_{Q_\theta} - k(x, \cdot), k_{Q_{\theta'}} - k(x', \cdot) \rangle_{\mathcal{H}} \leq 2C_k$.

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Proof.

- $h_{(\theta, \theta')}(x, x') := \langle k_{Q_{\theta}} - k(x, \cdot), k_{Q_{\theta'}} - k(x', \cdot) \rangle_{\mathcal{H}} \leq 2C_k$.

For any $\rho^2, \pi^2 \in \mathcal{M}_+^1(\Theta_{mult}^2)$, with probability $\geq 1 - \delta$

$$\text{kl}\left(\mathbb{E}_{(\theta, \theta') \sim \rho^2} \bar{U}_S(h_{(\theta, \theta')}) \middle\| \mathbb{E}_{(\theta, \theta') \sim \rho^2} \mathbb{E}_S \bar{U}_S(h_{(\theta, \theta')})\right) \leq \frac{2\text{KL}(\rho || \pi)}{\lfloor n/2 \rfloor} + c \quad (11)$$

where $c = \frac{\ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta}}{\lfloor n/2 \rfloor}$ and $\text{KL}(\rho^2 || \pi^2) = 2\text{KL}(\rho || \pi)$.

Is $\|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2$ a U -statistic?

Proof.

- $h_{(\theta, \theta')}(x, x') := \langle k_{Q_{\theta}} - k(x, \cdot), k_{Q_{\theta'}} - k(x', \cdot) \rangle_{\mathcal{H}} \leq 2C_k$.

For any $\rho^2, \pi^2 \in \mathcal{M}_+(\Theta_{mult}^2)$, with probability $\geq 1 - \delta$

$$\text{kl}\left(\mathbb{E}_{(\theta, \theta') \sim \rho^2} \bar{U}_S(h_{(\theta, \theta')}) \middle\| \mathbb{E}_{(\theta, \theta') \sim \rho^2} \mathbb{E}_S \bar{U}_S(h_{(\theta, \theta')})\right) \leq \frac{2\text{KL}(\rho \|\pi)}{\lfloor n/2 \rfloor} + c \quad (11)$$

where $c = \frac{\ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta}}{\lfloor n/2 \rfloor}$ and $\text{KL}(\rho^2 \|\pi^2) = 2\text{KL}(\rho \|\pi)$. Then

$$\frac{\mathbb{E}_{(\theta, \theta') \sim \rho^2}}{n(n-1)} \sum_{i \neq j}^n \langle k_{Q_{\theta}} - k(x_i, \cdot), k_{Q_{\theta'}} - k(x_j, \cdot) \rangle_{\mathcal{H}} = \frac{n \|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2 - V(\hat{\theta})}{(n-1)}$$

$$\frac{\mathbb{E}_{(\theta, \theta') \sim \rho^2}}{n(n-1)} \sum_{i \neq j}^n \mathbb{E}_S \langle k_{Q_{\theta}} - k(x_i, \cdot), k_{Q_{\theta'}} - k(x_j, \cdot) \rangle_{\mathcal{H}} = \|k_{Q_{\hat{\theta}}} - k_D\|_{\mathcal{H}}^2.$$

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching
- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2^{nd} Order U -Statistics
 - Do we have a U -statistic?
- 3 Close
 - **Choosing KL**
 - Conclusion

Choosing $\text{KL}(\hat{\rho}||\pi)$ for $\hat{\rho}, \pi \in \mathcal{M}_+^1(\Theta_{mult})$

Question

What are suitable choices for ρ and π such that $\text{KL}(\hat{\rho}||\pi)$ is manageable?

$$\text{kl}_{1/2} \left(\frac{n \|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2 - V(\hat{\theta})}{2C_k(n-1)} \left\| \frac{\|k_{Q_{\hat{\theta}}} - k_D\|_{\mathcal{H}}^2}{2C_k} \right. \right) \leq \frac{2\text{KL}(\hat{\rho}||\pi) + \ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta}}{\lfloor n/2 \rfloor}.$$

Choosing $\text{KL}(\hat{\rho}||\pi)$ for $\hat{\rho}, \pi \in \mathcal{M}_+^1(\Theta_{mult})$

Question

What are suitable choices for ρ and π such that $\text{KL}(\hat{\rho}||\pi)$ is manageable?

$$\text{kl}_{1/2} \left(\frac{n \|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2 - V(\hat{\theta})}{2C_k(n-1)} \left\| \frac{\|k_{Q_{\hat{\theta}}} - k_D\|_{\mathcal{H}}^2}{2C_k} \right. \right) \leq \frac{2\text{KL}(\hat{\rho}||\pi) + \ln \frac{\xi(\lfloor n/2 \rfloor)}{\delta}}{\lfloor n/2 \rfloor}.$$

Example (Dirichlet)

$\hat{\rho}, \pi \in \text{Dir}(\Theta_{mult})$ with respective means $\hat{\theta}$ and *uninformative* $\mathbf{1}/l$. Then

$$\text{KL}_{\text{Dir}}(\hat{\rho}||\pi) = \sum_{i=1}^l \log \frac{\Gamma(\hat{\theta}_i)}{\Gamma(1/l)} + \sum_{i=1}^l [\hat{\theta}_i - 1/l][\Psi(\hat{\theta}_i) - \Psi(1/l)], \quad (12)$$

where Γ and Ψ are the Gamma and Digamma functions.

Log-Normal Projection

- $\tau : \Theta_{mult} \rightarrow \mathbb{R}^l, \theta \mapsto (\ln(\frac{\theta_1}{g(\theta)}), \dots, \ln(\frac{\theta_l}{g(\theta)}))^\top, g(\theta) = (\theta_1 \dots \theta_l)^{1/l}.$

Log-Normal Projection

- $\tau : \Theta_{mult} \rightarrow \mathbb{R}^l, \theta \mapsto (\ln(\frac{\theta_1}{g(\theta)}), \dots, \ln(\frac{\theta_l}{g(\theta)}))^\top, g(\theta) = (\theta_1 \dots \theta_l)^{1/l}$.

Definition (Log-Normal Density)

Any $\rho \in LN(\Theta_{mult})$ with mean $\mu \in \Theta_{mult}$ and variance $1/\lambda$,

$$\frac{d\rho(\theta)}{d\theta} = |\nabla_{\theta}\tau| \left(\frac{2\pi}{\lambda}\right)^{-(l-1)/2} \exp\left(-\frac{\lambda}{2}\|\tau(\theta) - \tau(\mu)\|^2\right). \quad (13)$$

Log-Normal Projection

- $\tau : \Theta_{mult} \rightarrow \mathbb{R}^l, \theta \mapsto (\ln(\frac{\theta_1}{g(\theta)}), \dots, \ln(\frac{\theta_l}{g(\theta)}))^\top, g(\theta) = (\theta_1 \dots \theta_l)^{1/l}$.

Definition (Log-Normal Density)

Any $\rho \in LN(\Theta_{mult})$ with mean $\mu \in \Theta_{mult}$ and variance $1/\lambda$,

$$\frac{d\rho(\theta)}{d\theta} = |\nabla_{\theta}\tau| \left(\frac{2\pi}{\lambda}\right)^{-(l-1)/2} \exp\left(-\frac{\lambda}{2}\|\tau(\theta) - \tau(\mu)\|^2\right). \quad (13)$$

Example (Log-Normal KL)

$\hat{\rho}, \pi \in LN(\Theta_{mult})$ with respective means $\hat{\theta}$ and $\mathbf{1}/l$, and variance $1/\lambda$,

$$\text{KL}_{LN}(\hat{\rho}||\pi) = \frac{\lambda}{2}\|\tau(\hat{\theta})\|^2. \quad (14)$$

Outline

- 1 Motivation
 - Maximum Mean Discrepancy
 - Reproducing Kernel Moment Matching

- 2 First Order PAC-Bayes Bound
 - PAC-Bayes Bound for 2^{nd} Order U -Statistics
 - Do we have a U -statistic?

- 3 Close
 - Choosing KL
 - Conclusion

Conclusion

Summary:

- Derived a simple PAC-Bayes upper bound on the Maximum Mean Discrepancy between the input distribution and a mixture-based approximation.
- Given a class of prior and posterior with simple KL divergences.




Conclusion

Summary:

- Derived a simple PAC-Bayes upper bound on the Maximum Mean Discrepancy between the input distribution and a mixture-based approximation.
- Given a class of prior and posterior with simple KL divergences.

Future work:

- Find projections $\tau : \Theta_{mult} \rightarrow \mathbb{R}^l$ and distributions $\hat{\rho}$ and π that provide efficient $\text{KL}(\hat{\rho}||\pi)$ for optimisation.
- Examine bounding $\ln \mathbb{E}_S \exp \left(\lambda \left[\|k_{Q_{\hat{\theta}}} - k_D\|_{\mathcal{H}}^2 - \|k_{Q_{\hat{\theta}}} - k_{D_n}\|_{\mathcal{H}}^2 \right] \right)$.

-  Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B. & Smola, A. J. (2008).
CoRR abs/0805.2368.
-  Maurer, A. (2004).
-  Song, L., Zhang, X., Smola, A., Gretton, A. & Schölkopf, B. (2008).
In ICML '08: Proceedings of the 25th international conference on Machine learning pp. 992–999, ACM, New York, NY, USA.