PAC-Bayes Analysis: Links to Luckiness and Applications

John Shawe-Taylor University College London

PAC-Bayes Workshop March, 2010

Including joint work with Amiran Ambroladze and Emilio Parrado-Hernández, Cédric Archambeau, Matthew Higgs, Manfred Opper



- Definitions
- Priors from data distributions

2 Maximum entropy classification

- Generalisation
- Optimisation
- 3 GPs and SDEs
 - Gaussian Process regression
 - Variational approximation
 - Generalisation

Definitions Priors from data distributions

Luckiness definitions

Based on a function from samples and hypotheses:

 $L: X^m \times H \to \mathbb{R}^+,$

which measures the luckiness of a particular hypothesis with respect to the training examples.

• The level of luckiness is measured by seeing the number of functions that are luckier:

 $\ell(\mathbf{x},h) = |\{b \in \{0,1\}^m : \exists g \in H, g(\mathbf{x}) = b, L(\mathbf{x},g) \ge L(\mathbf{x},h)\}|.$

・ロト ・四ト ・ヨト ・ヨト

Definitions Priors from data distributions

Luckiness definitions

Based on a function from samples and hypotheses:

 $L: X^m \times H \to \mathbb{R}^+,$

which measures the luckiness of a particular hypothesis with respect to the training examples.

• The level of luckiness is measured by seeing the number of functions that are luckier:

 $\ell(\mathbf{x},h) = |\{b \in \{0,1\}^m : \exists g \in H, g(\mathbf{x}) = b, L(\mathbf{x},g) \ge L(\mathbf{x},h)\}|.$

・ロト ・四ト ・ヨト ・ヨト

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Example

- Motivating example was the case of large margin classifiers: luckiness measured margin on the sample
- Note that hyperplanes as classifiers cannot be ranked by margin until sample is seen
- Can overcome this difficulty if we consider real valued functions, require outputs to be ±1 and measure complexity by the norm of the weight vector
- This obscures the role of luckiness to capture alignment of hypotheses with the data generating distribution: eg density of distribution close to hyperplane
- Similar to idea of compatibility of Blum and Balcan (2005)
- Related to the local PAC-Bayes analysis of Catoni

Definitions Priors from data distributions

Defining priors from data distributions

- Can use part of the data to learn the prior: eg train svm on half the data and centre the prior gaussian on this weight vector (Emilio will give results for this)
- Can use some expectation over the true distribution to define the centre of the prior distribution such as

 $\mathbf{w} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[y\phi(\mathbf{x})]$

• or more sophisticated:

 $\mathbf{w} = \mathbb{E}_{S \sim \mathcal{D}^{m_0}}[\mathbf{w}_{\mathrm{SVM}}(S)]$

with $m_0 \ll m$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Definitions Priors from data distributions

Defining priors from data distributions

- Can use part of the data to learn the prior: eg train svm on half the data and centre the prior gaussian on this weight vector (Emilio will give results for this)
- Can use some expectation over the true distribution to define the centre of the prior distribution such as

 $\mathbf{W} = \mathbb{E}_{(\mathbf{X}, \mathbf{y}) \sim \mathcal{D}}[\mathbf{y}\phi(\mathbf{X})]$

• or more sophisticated:

 $\mathbf{w} = \mathbb{E}_{S \sim \mathcal{D}^{m_0}}[\mathbf{w}_{\mathrm{SVM}}(S)]$

with $m_0 \ll m$

Defining priors from data distributions

- Can use part of the data to learn the prior: eg train svm on half the data and centre the prior gaussian on this weight vector (Emilio will give results for this)
- Can use some expectation over the true distribution to define the centre of the prior distribution such as

$$\mathbf{w} = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\mathbf{y}\phi(\mathbf{x})]$$

or more sophisticated:

$$\mathbf{w} = \mathbb{E}_{S \sim \mathcal{D}^{m_0}}[\mathbf{w}_{\mathrm{SVM}}(S)]$$

with $m_0 \ll m$

A (1) > A (2) > A

Definitions Priors from data distributions

Defining priors from data distributions

In the latter cases use empirical versions for actual prior

- bound the difference between this and true (data distribution) prior
- use this to upper bound KL between true prior and posterior
- complexity term typically decays with increasing sample size
- Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size *m*
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size m
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size *m*
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size *m*
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size *m*
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size m
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size m
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Definitions Priors from data distributions

Defining priors from data distributions

- In the latter cases use empirical versions for actual prior
 - bound the difference between this and true (data distribution) prior
 - use this to upper bound KL between true prior and posterior
 - complexity term typically decays with increasing sample size
 - Shiliang will present theory and Emilio empirical results
- Can also define prior based on true risk (or expectation over SVM weight vectors on samples of training set size m
 - Even tighter bounds presented by Guy with extensions to manifold learning
 - Closer to Catoni's approach

Generalisation Optimisation

Maximum entropy learning

 $\bullet\,$ consider function class for ${\mathcal X}$ is a subset of the ℓ_∞ unit ball

$$\mathcal{F} = \left\{ f_{\boldsymbol{w}} : \boldsymbol{x} \in \mathcal{X} \mapsto \operatorname{sgn}\left(\sum_{i=1}^{N} w_i x_i\right) : \|\boldsymbol{w}\|_1 \leq 1 \right\},\$$

• want posterior distribution Q(w) such that can bound

 $P_{(\mathbf{x}, y) \sim \mathcal{D}}(f_{\mathbf{w}}(\mathbf{x}) \neq y) \leq 2e_{Q(\mathbf{w})}(= 2Q_{\mathcal{D}}(\mathbf{w})) = 2\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}, q \sim Q(\mathbf{w})}\left[I\left[q(\mathbf{x}) \neq y\right]\right]$

• Given a training sample $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we similarly define

$$\hat{e}_{Q(\boldsymbol{w})}(=\hat{Q}_{S}(\boldsymbol{w}))=\frac{1}{m}\sum_{i=1}^{m}\mathbb{E}_{q\sim Q(\boldsymbol{w})}\left[I\left[q(\boldsymbol{x}_{i})\neq y_{i}\right]\right].$$

Generalisation Optimisation

Maximum entropy learning

 $\bullet\,$ consider function class for ${\mathcal X}$ is a subset of the ℓ_∞ unit ball

$$\mathcal{F} = \left\{ f_{\boldsymbol{w}} : \boldsymbol{x} \in \mathcal{X} \mapsto \operatorname{sgn}\left(\sum_{i=1}^{N} w_i x_i\right) : \|\boldsymbol{w}\|_1 \leq 1 \right\},\$$

want posterior distribution Q(w) such that can bound

 $P_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}(f_{\boldsymbol{w}}(\boldsymbol{x})\neq\boldsymbol{y})\leq 2e_{Q(\boldsymbol{w})}(=2Q_{\mathcal{D}}(\boldsymbol{w}))=2\mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D},q\sim Q(\boldsymbol{w})}\left[I\left[q(\boldsymbol{x})\neq\boldsymbol{y}\right]\right]$

• Given a training sample $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we similarly define

$$\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})}(=\hat{Q}_{S}(\boldsymbol{w}))=\frac{1}{m}\sum_{i=1}^{m}\mathbb{E}_{q\sim Q(\boldsymbol{w})}\left[I\left[q(\boldsymbol{x}_{i})\neq y_{i}\right]\right].$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

Generalisation Optimisation

Maximum entropy learning

 $\bullet\,$ consider function class for ${\mathcal X}$ is a subset of the ℓ_∞ unit ball

$$\mathcal{F} = \left\{ f_{\boldsymbol{w}} : \boldsymbol{x} \in \mathcal{X} \mapsto \operatorname{sgn}\left(\sum_{i=1}^{N} w_i x_i\right) : \|\boldsymbol{w}\|_1 \leq 1 \right\},\$$

want posterior distribution Q(w) such that can bound

 $P_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}(f_{\boldsymbol{w}}(\boldsymbol{x})\neq\boldsymbol{y})\leq 2e_{Q(\boldsymbol{w})}(=2Q_{\mathcal{D}}(\boldsymbol{w}))=2\mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D},q\sim Q(\boldsymbol{w})}\left[I\left[q(\boldsymbol{x})\neq\boldsymbol{y}\right]\right]$

• Given a training sample $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, we similarly define

$$\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})}(=\hat{Q}_{S}(\boldsymbol{w}))=\frac{1}{m}\sum_{i=1}^{m}\mathbb{E}_{q\sim Q(\boldsymbol{w})}\left[I\left[q(\boldsymbol{x}_{i})\neq y_{i}\right]\right].$$

Generalisation Optimisation

Posterior distribution Q(w

• Classifier q involves random weight vector $W \in \mathbb{R}^N$ plus random threshold Θ

 $q_{W,\Theta}(\mathbf{x}) = \operatorname{sgn}\left(\langle W, \mathbf{x} \rangle - \Theta\right).$

• The distribution Q(w) of W will be discrete with

 $W = \operatorname{sgn}(w_i) e_i$; with probability $|w_i|, i = 1, \ldots, N$,

where e_i is the unit vector. The distribution of Θ is uniform on the interval [-1, 1].

Generalisation Optimisation

Posterior distribution Q(w

• Classifier *q* involves random weight vector $W \in \mathbb{R}^N$ plus random threshold Θ

 $q_{W,\Theta}(\mathbf{x}) = \operatorname{sgn}\left(\langle W, \mathbf{x} \rangle - \Theta\right).$

• The distribution Q(w) of W will be discrete with

 $W = \operatorname{sgn}(w_i)e_i$; with probability $|w_i|, i = 1, \ldots, N$,

where e_i is the unit vector. The distribution of Θ is uniform on the interval [-1, 1].

・ロット (母) ・ ヨ) ・ ・ ヨ)

Generalisation Optimisation

Error expression

Proposition

With the above definitions, we have for **w** satisfying $||w||_1 = 1$, that for any $(x, y) \in \mathcal{X} \times \{-1, +1\}$,

 $P_{q \sim Q(\boldsymbol{w})}(q(\boldsymbol{x}) \neq y) = 0.5(1 - y \langle \boldsymbol{w}, \boldsymbol{x} \rangle).$

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

<ロ> <同> <同> < 同> < 同> < 同> <

э

Generalisation Optimisation

Error expression proof

Proof.

$$P_{q \sim Q(\boldsymbol{w})}(q(\boldsymbol{x}) \neq \boldsymbol{y}) = \sum_{i=1}^{N} |w_i| P_{\Theta} \left(\operatorname{sgn} \left(\operatorname{sgn} (w_i) \langle \boldsymbol{e}_i, \boldsymbol{x} \rangle - \Theta \right) \neq \boldsymbol{y} \right) \\ = \sum_{i=1}^{N} |w_i| P_{\Theta} \left(\operatorname{sgn} \left(\operatorname{sgn} (w_i) x_i - \Theta \right) \neq \boldsymbol{y} \right) \\ = 0.5 \sum_{i=1}^{N} |w_i| (1 - y \operatorname{sgn} (w_i) x_i) \\ = 0.5 (1 - y \langle \boldsymbol{w}, \boldsymbol{x} \rangle),$$

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Generalisation Optimisation

Generalisation error

Corollary

$$P_{(\boldsymbol{x},y)\sim\mathcal{D}}\left(f_{\boldsymbol{w}}(\boldsymbol{x})\neq y
ight)\leq 2e_{Q(\boldsymbol{w})}.$$

Proof.

$$egin{aligned} & P_{q\sim Q(oldsymbol{w})}(q(oldsymbol{x})
eq y) &\geq 0.5 \ & \Leftrightarrow \ & f_{oldsymbol{w}}(oldsymbol{x}) &
eq y. \end{aligned}$$

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

크

Generalisation Optimisation

Base result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$\operatorname{KL}(\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})} \| \boldsymbol{e}_{Q(\boldsymbol{w})}) \leq \frac{\sum_{i=1}^{N} |w_i| \ln |w_i| + \ln(2N) + \ln((m+1)/\delta)}{m}$$

Proof.

- Use prior *P* uniform on unit vectors $\pm e_i$.
- Posterior described above so KL(P||Q(w)) equals ln(2N) entropy of w.

Generalisation Optimisation

Base result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$\operatorname{KL}(\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})} \| \boldsymbol{e}_{Q(\boldsymbol{w})}) \leq \frac{\sum_{i=1}^{N} |w_i| \ln |w_i| + \ln(2N) + \ln((m+1)/\delta)}{m}$$

Proof.

- Use prior *P* uniform on unit vectors $\pm e_i$.
- Posterior described above so KL(P||Q(w)) equals ln(2N)- entropy of w.

Generalisation Optimisation

Base result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$\operatorname{KL}(\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})} \| \boldsymbol{e}_{Q(\boldsymbol{w})}) \leq \frac{\sum_{i=1}^{N} |w_i| \ln |w_i| + \ln(2N) + \ln((m+1)/\delta)}{m}$$

Proof.

- Use prior *P* uniform on unit vectors $\pm e_i$.
- Posterior described above so KL(P||Q(w)) equals ln(2N) entropy of w.

Interpretation

- Suggests maximising the entropy as a means of minimising the bound.
- Problem that empirical error $\hat{e}_{Q(w)}$ is too large:

$$\hat{e}_{Q(\boldsymbol{w})} = \sum_{i=1}^{m} 0.5(1 - y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)$$

Generalisation

Optimisation

• Function of margin – but just linear function.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Interpretation

- Suggests maximising the entropy as a means of minimising the bound.
- Problem that empirical error $\hat{e}_{Q(w)}$ is too large:

$$\hat{e}_{Q(\boldsymbol{w})} = \sum_{i=1}^{m} 0.5(1 - y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)$$

Generalisation

Optimisation

• Function of margin – but just linear function.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Interpretation

- Suggests maximising the entropy as a means of minimising the bound.
- Problem that empirical error $\hat{e}_{Q(w)}$ is too large:

$$\hat{e}_{Q(\boldsymbol{w})} = \sum_{i=1}^{m} 0.5(1 - y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)$$

Generalisation

Optimisation

Function of margin – but just linear function.

Generalisation Optimisation

Boosting the bound

 Trick to boost the power of the bound is to take T independent samples of the distribution Q(w) and vote for the classification:

$$q_{\boldsymbol{W},\boldsymbol{\Theta}}(\boldsymbol{x}) = \operatorname{sgn}\left(\sum_{i=1}^{T}\operatorname{sgn}\left(\langle \boldsymbol{W}^{t}, \boldsymbol{x} \rangle - \boldsymbol{\Theta}^{t}\right)\right),$$

Now empirical error becomes

$$\hat{\boldsymbol{e}}_{\boldsymbol{Q}(\boldsymbol{w})} = \frac{0.5^{T}}{m} \sum_{i=1}^{m} \sum_{t=0}^{\lfloor T/2 \rfloor} {T \choose t} \left(1 + y_{i} \langle \boldsymbol{w}, \boldsymbol{x}_{i} \rangle\right)^{t} \left(1 - y_{i} \langle \boldsymbol{w}, \boldsymbol{x}_{i} \rangle\right)^{T-t},$$

giving sigmoid like loss as function of the margin.

(日) (四) (三) (三)
Generalisation Optimisation

Boosting the bound

 Trick to boost the power of the bound is to take T independent samples of the distribution Q(w) and vote for the classification:

$$q_{\boldsymbol{W},\boldsymbol{\Theta}}(\boldsymbol{x}) = \operatorname{sgn}\left(\sum_{i=1}^{T}\operatorname{sgn}\left(\langle \boldsymbol{W}^{t}, \boldsymbol{x} \rangle - \boldsymbol{\Theta}^{t}\right)\right),$$

Now empirical error becomes

$$\hat{\boldsymbol{e}}_{Q(\boldsymbol{w})} = \frac{0.5^{T}}{m} \sum_{i=1}^{m} \sum_{t=0}^{\lfloor T/2 \rfloor} {T \choose t} (1 + y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)^t (1 - y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)^{T-t},$$

giving sigmoid like loss as function of the margin.

< 同 > < 回 > < 回 >

Generalisation Optimisation

Full result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$P_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}(f_{\mathbf{w}}(\mathbf{x}) \neq \mathbf{y}) \leq 2KL^{-1}\left(\hat{e}_{Q^{T}(\mathbf{w})}, \frac{T\sum_{i=1}^{N} |w_{i}| \ln(|w_{i}|) + T\ln(2N) + \ln((m+1)/\delta)}{m}\right)$$

,

• Note penalty factor of T applied to KL

• Behaves like the (inverse) margin in usual bounds

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

(日)

Generalisation Optimisation

Full result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$P_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}\left(f_{\boldsymbol{w}}(\boldsymbol{x})\neq\boldsymbol{y}\right)\leq 2KL^{-1}\left(\hat{\boldsymbol{e}}_{\boldsymbol{Q}^{T}(\boldsymbol{w})},\frac{T\sum_{i=1}^{N}|w_{i}|\ln(|w_{i}|)+T\ln(2N)+\ln((m+1)/\delta)}{m}\right)$$

,

Note penalty factor of T applied to KL

• Behaves like the (inverse) margin in usual bounds

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Generalisation Optimisation

Full result

Theorem

With probability at least $1 - \delta$ over the draw of training sets of size *m*

$$P_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}}\left(f_{\boldsymbol{w}}(\boldsymbol{x})\neq\boldsymbol{y}\right)\leq 2KL^{-1}\left(\hat{\boldsymbol{e}}_{Q^{T}(\boldsymbol{w})},\frac{T\sum_{i=1}^{N}|\boldsymbol{w}_{i}|\ln(|\boldsymbol{w}_{i}|)+T\ln(2N)+\ln((m+1)/\delta)}{m}\right)$$

- Note penalty factor of T applied to KL
- Behaves like the (inverse) margin in usual bounds

• • • • • • • • • • • • •

Algorithmics

Bound motivates the optimisation:

 $\min_{\boldsymbol{w},\rho,\boldsymbol{\xi}} \qquad \sum_{j=1}^{N} |w_j| \ln |w_j| - C\rho + D \sum_{i=1}^{m} \xi_i$ subject to: $y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle \ge \rho - \xi_i, 1 \le i \le m,$ $\|\boldsymbol{w}\|_1 \le 1, \xi_i \ge 0, 1 \le i \le m.$

Generalisation

Optimisation

 This follows the SVM route of approximating the sigmoid like loss by the (convex) hinge loss

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Generalisation Optimisation

Algorithmics

Bound motivates the optimisation:

$$\min_{\boldsymbol{w},\rho,\boldsymbol{\xi}} \qquad \sum_{j=1}^{N} |w_j| \ln |w_j| - C\rho + D \sum_{i=1}^{m} \xi_i$$
subject to: $y_i \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle \ge \rho - \xi_i, 1 \le i \le m,$
 $\|\boldsymbol{w}\|_1 \le 1, \xi_i \ge 0, 1 \le i \le m.$

 This follows the SVM route of approximating the sigmoid like loss by the (convex) hinge loss

Generalisation Optimisation

Dual optimisation

$$\begin{split} \max_{\alpha} & L = -\sum_{j=1}^{N} \exp\left(\left|\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ij}\right| - 1 - \lambda\right) - \lambda\\ \text{subject to:} & \sum_{i=1}^{m} \alpha_{i} = C \quad 0 \leq \alpha_{i} \leq D, 1 \leq i \leq m. \end{split}$$

• Similar to SVM but with exponential function

- Surprisingly also gives dual sparsity
- Coordinate wise descent works very well (cf SMO algorithm)

A D N A D N A D N A D

Generalisation Optimisation

Dual optimisation

$$\begin{split} \max_{\alpha} & L = -\sum_{j=1}^{N} \exp\left(\left|\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ij}\right| - 1 - \lambda\right) - \lambda\\ \text{subject to:} & \sum_{i=1}^{m} \alpha_{i} = C \quad 0 \leq \alpha_{i} \leq D, 1 \leq i \leq m. \end{split}$$

Similar to SVM but with exponential function

- Surprisingly also gives dual sparsity
- Coordinate wise descent works very well (cf SMO algorithm)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Generalisation Optimisation

Dual optimisation

$$\begin{split} \max_{\alpha} & L = -\sum_{j=1}^{N} \exp\left(\left|\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ij}\right| - 1 - \lambda\right) - \lambda\\ \text{subject to:} & \sum_{i=1}^{m} \alpha_{i} = C \quad 0 \leq \alpha_{i} \leq D, 1 \leq i \leq m. \end{split}$$

- Similar to SVM but with exponential function
- Surprisingly also gives dual sparsity
- Coordinate wise descent works very well (cf SMO algorithm)

(日)

Generalisation Optimisation

Dual optimisation

$$\begin{split} \max_{\alpha} & L = -\sum_{j=1}^{N} \exp\left(\left|\sum_{i=1}^{m} \alpha_{i} y_{i} x_{ij}\right| - 1 - \lambda\right) - \lambda\\ \text{subject to:} & \sum_{i=1}^{m} \alpha_{i} = C \quad 0 \leq \alpha_{i} \leq D, 1 \leq i \leq m. \end{split}$$

- Similar to SVM but with exponential function
- Surprisingly also gives dual sparsity
- Coordinate wise descent works very well (cf SMO algorithm)

Generalisation Optimisation

Results: effect of varying



John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

< 冊

Generalisation Optimisation

Results

Bound and test errors:

Data	Bound	Error	SVM error	
Ionosphere	0.63	0.28	0.24	
Votes	0.78	0.35	0.35	
Glass	0.69	0.46	0.47	
Haberman	0.64	0.25	0.26	
Credit	0.60	0.25	0.28	

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Gaussian Process regression Variational approximation Generalisation

Gaussian Process Regression

- GP is distribution over real valued functions that is multivariate Gaussian when restricted to any finite subset of inputs
- Characterised by a kernel that specifies the covariance function when marginalising on any finite subset
- If have finite set of input/output observations generated with additive Gaussian noise on the outputs, posterior is also Gaussian process
- KL divergence between prior and posterior can be computed as (K = RR' is a Cholesky decomposition of K):

$$2\mathrm{KL}(\boldsymbol{Q}\|\boldsymbol{P}) = \log \det \left(\boldsymbol{I} + \frac{1}{\sigma^2}\boldsymbol{K}\right) - \mathrm{tr}\left(\left(\sigma^2\boldsymbol{I} + \boldsymbol{K}\right)^{-1}\boldsymbol{K}\right) + \left\|\boldsymbol{R}(\boldsymbol{K} + \sigma^2\boldsymbol{I})^{-1}\boldsymbol{y}\right\|^2$$

Gaussian Process regression Variational approximation Generalisation

Gaussian Process Regression

- GP is distribution over real valued functions that is multivariate Gaussian when restricted to any finite subset of inputs
- Characterised by a kernel that specifies the covariance function when marginalising on any finite subset
- If have finite set of input/output observations generated with additive Gaussian noise on the outputs, posterior is also Gaussian process
- KL divergence between prior and posterior can be computed as (K = RR' is a Cholesky decomposition of K):

$$2\mathrm{KL}(\boldsymbol{Q}\|\boldsymbol{P}) = \log \det \left(\boldsymbol{I} + \frac{1}{\sigma^2}\boldsymbol{K}\right) - \mathrm{tr}\left(\left(\sigma^2\boldsymbol{I} + \boldsymbol{K}\right)^{-1}\boldsymbol{K}\right) + \left\|\boldsymbol{R}(\boldsymbol{K} + \sigma^2\boldsymbol{I})^{-1}\boldsymbol{y}\right\|^2$$

Gaussian Process regression Variational approximation Generalisation

Gaussian Process Regression

- GP is distribution over real valued functions that is multivariate Gaussian when restricted to any finite subset of inputs
- Characterised by a kernel that specifies the covariance function when marginalising on any finite subset
- If have finite set of input/output observations generated with additive Gaussian noise on the outputs, posterior is also Gaussian process
- KL divergence between prior and posterior can be computed as (K = RR' is a Cholesky decomposition of K):

$$2\mathrm{KL}(Q||P) = \log \det \left(I + \frac{1}{\sigma^2} K \right) - \mathrm{tr} \left(\left(\sigma^2 I + K \right)^{-1} K \right) + \left\| R(K + \sigma^2 I)^{-1} \mathbf{y} \right\|^2$$

Gaussian Process regression Variational approximation Generalisation

Gaussian Process Regression

- GP is distribution over real valued functions that is multivariate Gaussian when restricted to any finite subset of inputs
- Characterised by a kernel that specifies the covariance function when marginalising on any finite subset
- If have finite set of input/output observations generated with additive Gaussian noise on the outputs, posterior is also Gaussian process
- KL divergence between prior and posterior can be computed as (K = RR' is a Cholesky decomposition of K):

$$2\mathrm{KL}(Q\|P) = \log\det\left(I + \frac{1}{\sigma^2}K\right) - \mathrm{tr}\left(\left(\sigma^2 I + K\right)^{-1}K\right) + \left\|R(K + \sigma^2 I)^{-1}\boldsymbol{y}\right\|^2$$

Gaussian Process regression Variational approximation Generalisation

Applying PAC-Bayes theorem

- Suggests can use the PB theorem if can create appropriate classifiers indexed by real value functions
- Consider for some $\epsilon > 0$ classifiers:

$$h_{f}^{\epsilon}(\boldsymbol{x}, \boldsymbol{y}) = \begin{cases} 1; & \text{if } |\boldsymbol{y} - f(\boldsymbol{x})| \leq \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

Can compute expected value of h^ε_f under posterior function:

$$\mathbb{E}_{f\sim Q}\left[h_f^{\epsilon}(\boldsymbol{x}, \boldsymbol{y})\right] = \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y} + \epsilon - \boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right) - \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y} - \epsilon - \boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right)$$

A D N A D N A D N A D

Gaussian Process regression Variational approximation Generalisation

Applying PAC-Bayes theorem

- Suggests can use the PB theorem if can create appropriate classifiers indexed by real value functions
- Consider for some $\epsilon > 0$ classifiers:

$$h_{f}^{\epsilon}(\boldsymbol{x}, \boldsymbol{y}) = \begin{cases} 1; & \text{if } |\boldsymbol{y} - f(\boldsymbol{x})| \leq \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

Can compute expected value of h^ε_f under posterior function:

$$\mathbb{E}_{f\sim Q}\left[h_f^{\epsilon}(\boldsymbol{x}, \boldsymbol{y})\right] = \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y} + \epsilon - \boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right) - \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y} - \epsilon - \boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right)$$

A D N A D N A D N A D

Gaussian Process regression Variational approximation Generalisation

Applying PAC-Bayes theorem

- Suggests can use the PB theorem if can create appropriate classifiers indexed by real value functions
- Consider for some $\epsilon > 0$ classifiers:

$$h_{f}^{\epsilon}(\boldsymbol{x}, \boldsymbol{y}) = \begin{cases} 1; & \text{if } |\boldsymbol{y} - f(\boldsymbol{x})| \leq \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

Can compute expected value of h^e_f under posterior function:

$$\mathbb{E}_{f\sim Q}\left[h_{f}^{\epsilon}(\boldsymbol{x},\boldsymbol{y})\right] = \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y}+\epsilon-\boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right) - \frac{1}{2} \operatorname{erf}\left(\frac{\boldsymbol{y}-\epsilon-\boldsymbol{m}(\boldsymbol{x})}{\sqrt{2\boldsymbol{v}(\boldsymbol{x})}}\right)$$

(4月) (1日) (1日)

Gaussian Process regression Variational approximation Generalisation

GP Result

Furthermore can lower bound expected value of point (*x*, *y*) in the posterior distribution by

$$2\epsilon \mathcal{N}(\boldsymbol{y}|\boldsymbol{m}(\boldsymbol{x}),\boldsymbol{v}(\boldsymbol{x})) \geq \mathbb{E}_{f\sim Q}\left[h_{f}^{\epsilon}(\boldsymbol{x},\boldsymbol{y})\right] - \sup_{\tau\in[\epsilon,\epsilon]} \frac{\epsilon^{2}}{2} \frac{2}{\boldsymbol{v}(\boldsymbol{x})\sqrt{2e\pi}}.$$

enabling an application of the PB Theorem to give:

$$\mathbb{E}\left[\mathcal{N}(\boldsymbol{y}|\boldsymbol{m}(\boldsymbol{x}),\boldsymbol{v}(\boldsymbol{x})) + \frac{\epsilon}{2\boldsymbol{v}(\boldsymbol{x})\sqrt{2\boldsymbol{e}\pi}}\right] \geq \frac{1}{2\epsilon} \mathrm{KL}^{-1}\left(\boldsymbol{E}(\epsilon),\frac{\boldsymbol{D} + \ln((\boldsymbol{m}+1)/\delta)}{\boldsymbol{m}}\right)$$

where $E(\epsilon)$ is the empirical average of $\mathbb{E}_{f\sim Q} \left[h_f^{\epsilon}(\mathbf{x}, y)\right]$ and D is the KL between prior and posterior.

(日)

Gaussian Process regression Variational approximation Generalisation

GP Experimental Results

- The robot arm problem (R), 150 training points and 51 test points.
- The Boston housing problem (H), 455 training points and 51 test points.
- The forest fire problem (F), 450 training points 67 test points.

Dat	σ	ê	KL^{-1}	<i>e</i> _{test}	KL ⁻¹	varGP etest
R	0.0494	0.8903	0.4782	0.8419		
Н	0.1924	0.8699	0.4645	0.7155	0.8401	0.9416
F	1.0129	0.5694	0.4557	0.5533		

Gaussian Process regression Variational approximation Generalisation

GP Experimental Results

 We can also plot the test accuracy and bound as a function of e:

Figure: Gaussian noise: Plot of $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[1 - \alpha(\mathbf{x})]$ against ϵ with for varying noise level η .



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Gaussian Process regression Variational approximation Generalisation

GP Experimental Results

• With Laplace noise:

Figure: Laplace noise: Plot of $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[1 - \alpha(\mathbf{x})]$ against ϵ with for varying η .



< □ > < □ > < □ > < □ > < □ >

Gaussian Process regression Variational approximation Generalisation

GP Experimental Results

• Robot arm problem and Boston Housing:

Figure: Confidence levels for Robot arm problem



Gaussian Process regression Variational approximation Generalisation

Stochastic Differential Equation Models

• Consider modelling a time varying process with a (non-linear) stochastic differential equation:

 $d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + \sqrt{\mathbf{\Sigma}} \ d\mathbf{W}$

f(x, t) is a non-linear drift term and dW is a Wiener process
This is the limit of the discrete time equation:

 $\Delta \mathbf{x}_k \equiv \mathbf{x}_{k+1} - \mathbf{x}_k = \mathbf{f}(\mathbf{x}_k) \Delta t + \sqrt{\Delta t \, \mathbf{\Sigma}} \, \boldsymbol{\epsilon}_k \; .$

where ϵ_k is zero mean, unit variance Gaussian noise.

(日) (四) (日) (日) (日)

Gaussian Process regression Variational approximation Generalisation

Stochastic Differential Equation Models

• Consider modelling a time varying process with a (non-linear) stochastic differential equation:

 $d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + \sqrt{\mathbf{\Sigma}} d\mathbf{W}$

f(x, t) is a non-linear drift term and dW is a Wiener process
This is the limit of the discrete time equation:

 $\Delta \mathbf{x}_k \equiv \mathbf{x}_{k+1} - \mathbf{x}_k = \mathbf{f}(\mathbf{x}_k) \Delta t + \sqrt{\Delta t \, \mathbf{\Sigma}} \, \boldsymbol{\epsilon}_k \; .$

where ϵ_k is zero mean, unit variance Gaussian noise.

(日) (四) (日) (日) (日)

Gaussian Process regression Variational approximation Generalisation

Stochastic Differential Equation Models

• Consider modelling a time varying process with a (non-linear) stochastic differential equation:

 $d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + \sqrt{\mathbf{\Sigma}} d\mathbf{W}$

- **f**(**x**, *t*) is a non-linear drift term and *d***W** is a Wiener process
- This is the limit of the discrete time equation:

$$\Delta \mathbf{x}_k \equiv \mathbf{x}_{k+1} - \mathbf{x}_k = \mathbf{f}(\mathbf{x}_k) \Delta t + \sqrt{\Delta t \, \boldsymbol{\Sigma}} \, \boldsymbol{\epsilon}_k \, .$$

where ϵ_k is zero mean, unit variance Gaussian noise.

Gaussian Process regression Variational approximation Generalisation

Variational approximation

 We use the Bayesian approach to data modelling with a noise model given by:

 $\rho(\mathbf{y}_n|\mathbf{x}(t_n)) = \mathcal{N}(\mathbf{y}_n|\mathbf{H}\mathbf{x}(t_n),\mathbf{R}),$

• We consider a variational approximation of the posterior using a time-varying linear SDE:

 $d\mathbf{x} = \mathbf{f}_L(\mathbf{x}, t) dt + \sqrt{\mathbf{\Sigma}} d\mathbf{W},$

where

$$\mathbf{f}_L(\mathbf{x},t) = -\mathbf{A}(t)\mathbf{x} + \mathbf{b}(t).$$

Gaussian Process regression Variational approximation Generalisation

Variational approximation

 We use the Bayesian approach to data modelling with a noise model given by:

```
\rho(\mathbf{y}_n|\mathbf{x}(t_n)) = \mathcal{N}(\mathbf{y}_n|\mathbf{H}\mathbf{x}(t_n),\mathbf{R}),
```

• We consider a variational approximation of the posterior using a time-varying linear SDE:

$$d\mathbf{x} = \mathbf{f}_L(\mathbf{x}, t) dt + \sqrt{\mathbf{\Sigma}} \ d\mathbf{W},$$

where

$$\mathbf{f}_L(\mathbf{x},t) = -\mathbf{A}(t)\mathbf{x} + \mathbf{b}(t).$$

Gaussian Process regression Variational approximation Generalisation

Girsanov change of measure

- Measure for the drift f denoted by P and the one for drift f_L by Q.
- The KL divergence in the infinite dimensional setting is given by Radon-Nikodym derivative of Q with respect to P:

$$\operatorname{KL}[Q \| P] = \int dQ \ln \frac{dQ}{dP} = E_Q \ln \frac{dQ}{dP} ,$$

which can be computed as

$$\frac{dQ}{dP} = \exp\left\{-\int_{t_0}^{t_f} (\mathbf{f} - \mathbf{f}_L)^\top \mathbf{\Sigma}^{-1/2} \ d\widehat{W}_t + \frac{1}{2}\int_{t_0}^{t_f} (\mathbf{f} - \mathbf{f}_L)^\top \mathbf{\Sigma}^{-1} (\mathbf{f} - \mathbf{f}_L) \ dt\right\},$$

where \widehat{W} is a Wiener process with respect to Q.

Gaussian Process regression Variational approximation Generalisation

Girsanov change of measure

- Measure for the drift f denoted by P and the one for drift f_L by Q.
- The KL divergence in the infinite dimensional setting is given by Radon-Nikodym derivative of Q with respect to P:

$$\mathrm{KL}[\boldsymbol{Q}\|\boldsymbol{P}] = \int d\boldsymbol{Q} \ln \frac{d\boldsymbol{Q}}{d\boldsymbol{P}} = \boldsymbol{E}_{\boldsymbol{Q}} \ln \frac{d\boldsymbol{Q}}{d\boldsymbol{P}} ,$$

which can be computed as

$$\frac{dQ}{dP} = \exp\left\{-\int_{t_0}^{t_f} (\mathbf{f} - \mathbf{f}_L)^\top \mathbf{\Sigma}^{-1/2} \ d\widehat{W}_t + \frac{1}{2}\int_{t_0}^{t_f} (\mathbf{f} - \mathbf{f}_L)^\top \mathbf{\Sigma}^{-1} (\mathbf{f} - \mathbf{f}_L) \ dt\right\},$$

where \widehat{W} is a Wiener process with respect to Q.

A B > A B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A
 B > A

Gaussian Process regression Variational approximation Generalisation

KL divergence

Hence, KL divergence is

$\mathrm{KL}[\boldsymbol{Q}\|\boldsymbol{P}] = \frac{1}{2} \int_{t_0}^{t_f} \left\langle (\mathbf{f}(\mathbf{x}(t), t) - \mathbf{f}_L(\mathbf{x}(t), t))^\top \boldsymbol{\Sigma}^{-1} (\mathbf{f}(\mathbf{x}(t), t) - \mathbf{f}_L(\mathbf{x}(t), t)) \right\rangle_{a_t}$

where $\langle \cdot \rangle_{q_t}$ denotes the expectation with respect to the marginal density at time *t* of the measure *Q*.

Gaussian Process regression Variational approximation Generalisation

Variational approximation

 As approximating SDE is linear, marginal distribution *q_t* is Gaussian

 $q_t(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}(t), \mathbf{S}(t)).$

 with the mean m(t) and covariance S(t) described by ordinary differential equations (ODEs):

$$\frac{d\mathbf{m}}{dt} = -\mathbf{A}\mathbf{m} + \mathbf{b},$$

$$\frac{d\mathbf{S}}{dt} = -\mathbf{A}\mathbf{S} - \mathbf{S}\mathbf{A}^{\mathrm{T}} + \boldsymbol{\Sigma}$$

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Gaussian Process regression Variational approximation Generalisation

Variational approximation

 As approximating SDE is linear, marginal distribution *q_t* is Gaussian

 $q_t(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}(t), \mathbf{S}(t)).$

 with the mean m(t) and covariance S(t) described by ordinary differential equations (ODEs):

$$\frac{d\mathbf{m}}{dt} = -\mathbf{A}\mathbf{m} + \mathbf{b},$$

$$\frac{d\mathbf{S}}{dt} = -\mathbf{A}\mathbf{S} - \mathbf{S}\mathbf{A}^{\mathrm{T}} + \mathbf{\Sigma}$$

Gaussian Process regression Variational approximation Generalisation

Algorithmics

- Using Lagrangian methods can derive algorithm that finds the variational approximation by minimising the KL divergence between posterior and approximating distribution.
- But KL also appears in the PAC-Bayes bound is it possible to define appropriate loss over paths ω that captures the properties of interest?

Gaussian Process regression Variational approximation Generalisation

Algorithmics

- Using Lagrangian methods can derive algorithm that finds the variational approximation by minimising the KL divergence between posterior and approximating distribution.
- But KL also appears in the PAC-Bayes bound is it possible to define appropriate loss over paths ω that captures the properties of interest?
Error estimation

For ω : [0, T] → ℝ^D defining a trajectory ω(t) ∈ ℝ^D, we define the classifier h_ω by

$$h_{\omega}(\mathbf{y}, t) = \begin{cases} 1; & \text{if } \|\mathbf{y} - \mathbf{H}_{\omega}(t)\| \leq \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

- where the actual observations are linear functions of the state variable given by the operator **H**.
- Prior and posterior distribution over functions are inherited from distributions P and Q over paths ω .
- Hence, $P = p_{sde}$ and Q = q defined by linear approximating sde.

A D N A D N A D N A D

Error estimation

For ω : [0, T] → ℝ^D defining a trajectory ω(t) ∈ ℝ^D, we define the classifier h_ω by

$$h_{\omega}(\mathbf{y}, t) = \begin{cases} 1; & \text{if } \|\mathbf{y} - \mathbf{H}_{\omega}(t)\| \leq \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

- where the actual observations are linear functions of the state variable given by the operator H.
- Prior and posterior distribution over functions are inherited from distributions P and Q over paths ω .
- Hence, P = p_{sde} and Q = q defined by linear approximating sde.

(日) (四) (三) (三)

Error estimation

For ω : [0, T] → ℝ^D defining a trajectory ω(t) ∈ ℝ^D, we define the classifier h_ω by

$$h_{\omega}(\mathbf{y}, t) = \begin{cases} 1; & \text{if } \|\mathbf{y} - \mathbf{H}_{\omega}(t)\| \le \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

- where the actual observations are linear functions of the state variable given by the operator H.
- Prior and posterior distribution over functions are inherited from distributions *P* and *Q* over paths ω.
- Hence, P = p_{sde} and Q = q defined by linear approximating sde.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Error estimation

For ω : [0, T] → ℝ^D defining a trajectory ω(t) ∈ ℝ^D, we define the classifier h_ω by

$$h_{\omega}(\mathbf{y}, t) = \begin{cases} 1; & \text{if } \|\mathbf{y} - \mathbf{H}_{\omega}(t)\| \le \epsilon; \\ 0; & \text{otherwise.} \end{cases}$$

- where the actual observations are linear functions of the state variable given by the operator H.
- Prior and posterior distribution over functions are inherited from distributions *P* and *Q* over paths ω.
- Hence, P = p_{sde} and Q = q defined by linear approximating sde.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Gaussian Process regression Variational approximation Generalisation

Generalisation analysis

• For the PAC-Bayes analysis we must compute: KL(Q||P), e_Q , \hat{e}_Q . We have as above

$$\mathrm{KL}(Q\|P) = \int dq \ln \frac{dq}{dp_{\mathrm{sde}}}$$

• If we now consider a fixed sample (y, t) we can estimate

$$\mathbb{E}_{\omega \sim \mathcal{Q}}\left[h_{\omega}(\mathbf{y},t)
ight] = \int I\left[\left\|\mathbf{H}\mathbf{x}-\mathbf{y}
ight\| \leq \epsilon
ight] dq_t(\mathbf{x}),$$

• For sufficiently small values of ϵ we can approximate by

 $\approx \frac{V_d \epsilon^d}{(2\pi |\mathsf{HS}(t)\mathsf{H}^{\mathsf{T}}|)^{d/2}} \exp\left(-(\mathbf{y} - \mathsf{Hm}(t))^{\mathsf{T}} (\mathsf{HS}(t)\mathsf{H}^{\mathsf{T}})^{-1} (\mathbf{y} - \mathsf{Hm}(t))\right),$

where V_d is the volume of a unit ball in \mathbb{R}^d .

Gaussian Process regression Variational approximation Generalisation

Generalisation analysis

• For the PAC-Bayes analysis we must compute: KL(Q||P), e_Q , \hat{e}_Q . We have as above

$$\mathrm{KL}(Q\|P) = \int dq \ln \frac{dq}{dp_{\mathrm{sde}}}$$

• If we now consider a fixed sample (y, t) we can estimate

$$\mathbb{E}_{\omega \sim \mathcal{Q}}\left[h_{\omega}(\mathbf{y},t)
ight] = \int I\left[\left\|\mathbf{H}\mathbf{x}-\mathbf{y}
ight\| \leq \epsilon
ight] dq_t(\mathbf{x}),$$

• For sufficiently small values of ϵ we can approximate by

 $\approx \frac{V_d \epsilon^d}{(2\pi |\mathsf{HS}(t)\mathsf{H}^{\mathsf{T}}|)^{d/2}} \exp\left(-(\mathbf{y} - \mathsf{Hm}(t))^{\mathsf{T}} (\mathsf{HS}(t)\mathsf{H}^{\mathsf{T}})^{-1} (\mathbf{y} - \mathsf{Hm}(t))\right),$

where V_d is the volume of a unit ball in \mathbb{R}^d .

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

Gaussian Process regression Variational approximation Generalisation

Generalisation analysis

• For the PAC-Bayes analysis we must compute: KL(Q||P), e_Q , \hat{e}_Q . We have as above

$$\mathrm{KL}(Q\|P) = \int dq \ln \frac{dq}{dp_{\mathrm{sde}}}$$

• If we now consider a fixed sample (y, t) we can estimate

$$\mathbb{E}_{\omega \sim Q}\left[h_{\omega}(\mathbf{y}, t)\right] = \int I\left[\|\mathbf{H}\mathbf{x} - \mathbf{y}\| \le \epsilon\right] dq_t(\mathbf{x}),$$

• For sufficiently small values of ϵ we can approximate by

 $\approx \frac{V_d \epsilon^d}{(2\pi |\mathbf{HS}(t)\mathbf{H}^T|)^{d/2}} \exp\left(-(\mathbf{y} - \mathbf{Hm}(t))^T (\mathbf{HS}(t)\mathbf{H}^T)^{-1} (\mathbf{y} - \mathbf{Hm}(t))\right),$

where V_d is the volume of a unit ball in \mathbb{R}^d .

Gaussian Process regression Variational approximation Generalisation

Error estimates

Note that e_Q is simply

 $\boldsymbol{e}_{Q} = \mathbb{E}_{(\mathbf{y},t)\sim\mu}\mathbb{E}_{\omega\sim Q}\left[\boldsymbol{h}_{\omega}(\mathbf{y},t)\right] \propto \int \mathcal{N}(\mathbf{y}|\mathbf{Hm}(t),\mathbf{HS}(t)\mathbf{H}^{T})d\mu(\mathbf{y},t),$

while \hat{e}_Q is the empirical average of this quantity.

- A tension arises in setting *ϵ* − if large approximation inaccurate.
- If e_Q and \hat{e}_Q both small, the bound implied by the $\operatorname{KL}(e_Q \| \hat{e}_Q) \leq C$ becomes weak.

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・ ・

Gaussian Process regression Variational approximation Generalisation

Error estimates

Note that e_Q is simply

 $\boldsymbol{e}_{\boldsymbol{Q}} = \mathbb{E}_{(\mathbf{y},t)\sim\mu}\mathbb{E}_{\omega\sim\boldsymbol{Q}}\left[\boldsymbol{h}_{\omega}(\mathbf{y},t)\right] \propto \int \mathcal{N}(\mathbf{y}|\mathbf{Hm}(t),\mathbf{HS}(t)\mathbf{H}^{T})d\mu(\mathbf{y},t),$

while \hat{e}_Q is the empirical average of this quantity.

- A tension arises in setting *ϵ* − if large approximation inaccurate.
- If e_Q and \hat{e}_Q both small, the bound implied by the $\operatorname{KL}(e_Q \| \hat{e}_Q) \leq C$ becomes weak.

・ロ・ ・ 四・ ・ 回・ ・ 回・

Gaussian Process regression Variational approximation Generalisation

Error estimates

Note that e_Q is simply

 $\boldsymbol{e}_{\boldsymbol{Q}} = \mathbb{E}_{(\mathbf{y},t)\sim\mu}\mathbb{E}_{\omega\sim\boldsymbol{Q}}\left[\boldsymbol{h}_{\omega}(\mathbf{y},t)\right] \propto \int \mathcal{N}(\mathbf{y}|\mathbf{Hm}(t),\mathbf{HS}(t)\mathbf{H}^{T})d\mu(\mathbf{y},t),$

while \hat{e}_Q is the empirical average of this quantity.

- A tension arises in setting *ϵ* if large approximation inaccurate.
- If e_Q and \hat{e}_Q both small, the bound implied by the $KL(e_Q \| \hat{e}_Q) \le C$ becomes weak.

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・ ・

Gaussian Process regression Variational approximation Generalisation

Refining the distributions

- Overcome this weakness by taking *K*-fold product distributions and defining h_(ω1,...,ωκ) as
- $h_{(\omega_1,...,\omega_K)}(\mathbf{y},t) = \begin{cases} 1; & \text{if there exists } 1 \le i \le K \text{ such that} \|\mathbf{y} \mathbf{H}\omega_i(t)\| \le \epsilon; \\ 0; & \text{otherwise.} \end{cases}$

• We now have

 $\mathbb{E}_{(\omega_1,...,\omega_K)\sim Q^K} \left[h_{(\omega_1,...,\omega_K)}(\mathbf{y},t) \right] \approx 1 - \left(1 - \int I[\|\mathbf{H}\mathbf{x} - \mathbf{y}\| \le \epsilon] \, dq_t(\mathbf{x}) \right)^K \\ \approx K V_d \epsilon^d \mathcal{N}(\mathbf{y} | \mathbf{Hm}(t), \mathbf{HS}(t) \mathbf{H}^T),$

(日)

Gaussian Process regression Variational approximation Generalisation

Refining the distributions

Overcome this weakness by taking *K*-fold product distributions and defining h_{(ω1,...,ωK}) as

 $h_{(\omega_1,...,\omega_K)}(\mathbf{y},t) = \begin{cases} 1; & \text{if there exists } 1 \le i \le K \text{ such that} \|\mathbf{y} - \mathbf{H}\omega_i(t)\| \le \epsilon; \\ 0; & \text{otherwise.} \end{cases}$

• We now have

 $\mathbb{E}_{(\omega_1,...,\omega_K)\sim Q^K} \left[h_{(\omega_1,...,\omega_K)}(\mathbf{y},t) \right] \approx 1 - \left(1 - \int I[\|\mathbf{H}\mathbf{x} - \mathbf{y}\| \le \epsilon] \, dq_t(\mathbf{x}) \right)^K \\ \approx K V_d \epsilon^d \mathcal{N}(\mathbf{y} | \mathbf{Hm}(t), \mathbf{HS}(t) \mathbf{H}^T),$

ヘロト 人間 ト 人造 ト 人造 ト

Final result

• Putting all together gives final bound:

$$\begin{split} \mathbb{E}_{(\mathbf{y},t)\sim\mu} \left[\mathcal{N}(\mathbf{y}|\mathbf{Hm}(t),\mathbf{HS}(t)\mathbf{H}^{T}) \right] \geq & \\ \frac{1}{V_{d}\epsilon^{d}K} \mathrm{KL}^{-1} \left(\mathcal{K}V_{d}\epsilon^{d} \hat{\mathbb{E}} \left[\mathcal{N}(\mathbf{y}|\mathbf{Hm}(t),\mathbf{HS}(t)\mathbf{H}^{T}) \right], \\ & \frac{\mathcal{K} \int_{0}^{T} \mathcal{E}_{\mathrm{sde}}(t) dt + \ln((m+1)/\delta)}{m} \right). \end{split}$$

where

$$E_{sde}(t) = \frac{1}{2} \left\langle (\mathbf{f}(\mathbf{x}) - \mathbf{f}_{L}(\mathbf{x}, t))^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{f}(\mathbf{x}) - \mathbf{f}_{L}(\mathbf{x}, t)) \right\rangle_{q_{t}},$$

John Shawe-Taylor University College London PAC-Bayes Analysis: Links to Luckiness and Applications

・ロ・ ・ 四・ ・ 回・ ・ 回・

Gaussian Process regression Variational approximation Generalisation

Small scale experiment

• We applied the analysis to the results of performing a variational Bayesian approximation to the Lorentz attractor in three dimension. The quality of the fit with 49 examples was good.



PAC-Bayes Analysis: Links to Luckiness and Applications

Gaussian Process regression Variational approximation Generalisation

Small scale experiment

- chose V_d e^d to optimise the bound fairly small ball implying that our approximation should be reasonable.
- compared the bound with the left hand side estimated on a random draw of 99 test points. The corresponding values are

т	dt	ê _Q	Α	e _Q	$\mathrm{KL}^{-1}(\cdot,\cdot)/V$
49	0.005	0.137	3.536	0.128	0.004

Gaussian Process regression Variational approximation Generalisation

Small scale experiment

- chose V_d e^d to optimise the bound fairly small ball implying that our approximation should be reasonable.
- compared the bound with the left hand side estimated on a random draw of 99 test points. The corresponding values are

т	dt	ê _Q	A	e _Q	$\mathrm{KL}^{-1}(\cdot,\cdot)/V$
49	0.005	0.137	3.536	0.128	0.004

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE – hence Gaussian process posterior.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE – hence Gaussian process posterior.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE – hence Gaussian process posterior.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE – hence Gaussian process posterior.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE hence Gaussian process posterior.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Gaussian Process regression Variational approximation Generalisation

Conclusions

- Links between luckiness and choosing the prior based on the data distribution
- Applications to maximum entropy classification
- Also consider lower bounding the accuracy of a posterior distribution for Gaussian processes (GP)
- Applied the theory to bound the performance of estimations made using approximate Bayesian inference for dynamical systems:
 - Prior determined by a non-linear stochastic differential equation (SDE)
 - Variational approximation results in posterior given by an approximating linear SDE – hence Gaussian process posterior.