Data Dependent Priors in PAC-Bayes Bounds

John Shawe-Taylor University College London

Joint work with Emilio Parrado-Hernández and Amiran Ambroladze

August, 2010

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

A (10) + A (10) +



- 2 PAC-Bayes Analysis
 - Definitions
 - PAC-Bayes Theorem
 - Proof outline
 - Applications



Linear Classifiers

- General Approach
- Learning the prior
- New prior for linear functions
- Prior-SVM

3 1 4 3 1

Evidence and generalisation

- Link between evidence and generalisation hypothesised by McKay
- First formal link was obtained by S-T & Williamson (1997): PAC Analysis of a Bayes Estimator
- Bound on generalisation in terms of the volume of the sphere that can be inscribed in the version space – included a dependence on the dimensionality of the space
- Used Luckiness framework a data-dependent style of frequentist bound also used to bound generalisation of SVMs for which no dependence on the dimensionality is needed, just on the margin

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Evidence and generalisation

- Link between evidence and generalisation hypothesised by McKay
- First formal link was obtained by S-T & Williamson (1997): PAC Analysis of a Bayes Estimator
- Bound on generalisation in terms of the volume of the sphere that can be inscribed in the version space – included a dependence on the dimensionality of the space
- Used Luckiness framework a data-dependent style of frequentist bound also used to bound generalisation of SVMs for which no dependence on the dimensionality is needed, just on the margin

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Evidence and generalisation

- Link between evidence and generalisation hypothesised by McKay
- First formal link was obtained by S-T & Williamson (1997): PAC Analysis of a Bayes Estimator
- Bound on generalisation in terms of the volume of the sphere that can be inscribed in the version space – included a dependence on the dimensionality of the space
- Used Luckiness framework a data-dependent style of frequentist bound also used to bound generalisation of SVMs for which no dependence on the dimensionality is needed, just on the margin

(日本) (日本) (日本)

Evidence and generalisation

- Link between evidence and generalisation hypothesised by McKay
- First formal link was obtained by S-T & Williamson (1997): PAC Analysis of a Bayes Estimator
- Bound on generalisation in terms of the volume of the sphere that can be inscribed in the version space – included a dependence on the dimensionality of the space
- Used Luckiness framework a data-dependent style of frequentist bound also used to bound generalisation of SVMs for which no dependence on the dimensionality is needed, just on the margin

・ 同 ト ・ ヨ ト ・ ヨ ト ・

Definitions PAC-Bayes Theorem Proof outline Applications

PAC-Bayes Theorem

First version proved by McAllester in 1999

- Improved proof and bound due to Seeger in 2002 with application to Gaussian processes
- Application to SVMs by Langford and S-T also in 2002
- Excellent tutorial by Langford appeared in 2005 in JMLR

Definitions PAC-Bayes Theorem Proof outline Applications

PAC-Bayes Theorem

- First version proved by McAllester in 1999
- Improved proof and bound due to Seeger in 2002 with application to Gaussian processes
- Application to SVMs by Langford and S-T also in 2002
- Excellent tutorial by Langford appeared in 2005 in JMLR

Definitions PAC-Bayes Theorem Proof outline Applications

PAC-Bayes Theorem

- First version proved by McAllester in 1999
- Improved proof and bound due to Seeger in 2002 with application to Gaussian processes
- Application to SVMs by Langford and S-T also in 2002
- Excellent tutorial by Langford appeared in 2005 in JMLR

Definitions PAC-Bayes Theorem Proof outline Applications

PAC-Bayes Theorem

- First version proved by McAllester in 1999
- Improved proof and bound due to Seeger in 2002 with application to Gaussian processes
- Application to SVMs by Langford and S-T also in 2002
- Excellent tutorial by Langford appeared in 2005 in JMLR

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Prior and posterior distributions

- The PAC-Bayes theorem involves a class of classifiers *C* together with a prior distribution *P* and posterior *Q* over *C*
- The distribution *P* must be chosen before learning, but the bound holds for all choices of *Q*, hence *Q* does not need to be the classical Bayesian posterior
- The bound holds for all (prior) choices of P hence it's validity is not affected by a poor choice of P though the quality of the resulting bound may be

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Prior and posterior distributions

- The PAC-Bayes theorem involves a class of classifiers *C* together with a prior distribution *P* and posterior *Q* over *C*
- The distribution *P* must be chosen before learning, but the bound holds for all choices of *Q*, hence *Q* does not need to be the classical Bayesian posterior
- The bound holds for all (prior) choices of P hence it's validity is not affected by a poor choice of P though the quality of the resulting bound may be

• □ > • □ > • □ > • □ > • □ > ·

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Prior and posterior distributions

- The PAC-Bayes theorem involves a class of classifiers *C* together with a prior distribution *P* and posterior *Q* over *C*
- The distribution *P* must be chosen before learning, but the bound holds for all choices of *Q*, hence *Q* does not need to be the classical Bayesian posterior
- The bound holds for all (prior) choices of P hence it's validity is not affected by a poor choice of P though the quality of the resulting bound may be

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Error measures

- Being a frequentist (PAC) style result we assume an unknown distribution \mathcal{D} on the input space *X*.
- \mathcal{D} is used to generate the labelled training samples i.i.d., i.e. $S \sim \mathcal{D}^m$
- It is also used to measure generalisation error c_D of a classifier c:

$$c_{\mathcal{D}} = \Pr_{(x,y) \sim \mathcal{D}}(c(x) \neq y)$$

• The empirical generalisation error is denoted \hat{c}_S :

$$\hat{c}_S = \frac{1}{m} \sum_{(x,y)\in S} I[c(x) \neq y]$$
 where $I[\cdot]$ indicator function.

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Error measures

- Being a frequentist (PAC) style result we assume an unknown distribution \mathcal{D} on the input space X.
- \mathcal{D} is used to generate the labelled training samples i.i.d., i.e. $S \sim \mathcal{D}^m$
- It is also used to measure generalisation error c_D of a classifier c:

$$c_{\mathcal{D}} = \Pr_{(x,y)\sim\mathcal{D}}(c(x) \neq y)$$

• The empirical generalisation error is denoted \hat{c}_S :

$$\hat{c}_S = \frac{1}{m} \sum_{(x,y) \in S} I[c(x) \neq y]$$
 where $I[\cdot]$ indicator function.

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Error measures

- Being a frequentist (PAC) style result we assume an unknown distribution \mathcal{D} on the input space X.
- \mathcal{D} is used to generate the labelled training samples i.i.d., i.e. $S \sim \mathcal{D}^m$
- It is also used to measure generalisation error c_D of a classifier c:

$$c_{\mathcal{D}} = \Pr_{(x,y)\sim\mathcal{D}}(c(x) \neq y)$$

• The empirical generalisation error is denoted \hat{c}_S :

 $\hat{c}_S = \frac{1}{m} \sum_{(x,y) \in S} I[c(x) \neq y]$ where $I[\cdot]$ indicator function.

< ロ > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Error measures

- Being a frequentist (PAC) style result we assume an unknown distribution \mathcal{D} on the input space X.
- \mathcal{D} is used to generate the labelled training samples i.i.d., i.e. $S \sim \mathcal{D}^m$
- It is also used to measure generalisation error c_D of a classifier c:

$$c_{\mathcal{D}} = \Pr_{(x,y) \sim \mathcal{D}}(c(x) \neq y)$$

• The empirical generalisation error is denoted \hat{c}_S :

$$\hat{c}_{S} = \frac{1}{m} \sum_{(x,y)\in S} I[c(x) \neq y]$$
 where $I[\cdot]$ indicator function.

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Assessing the posterior

- The result is concerned with bounding the performance of a probabilistic classifier that given a test input *x* chooses a classifier *c* ~ *Q* (the posterior) and returns *c*(*x*)
- We are interested in the relation between two quantities:

$$Q_{\mathcal{D}} = \mathbb{E}_{c \sim Q}[c_{\mathcal{D}}]$$

the true error rate of the probabilistic classifier and

$$\hat{Q}_{\mathcal{S}} = \mathbb{E}_{c \sim Q}[\hat{c}_{\mathcal{S}}]$$

(日)

its empirical error rate

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Assessing the posterior

- The result is concerned with bounding the performance of a probabilistic classifier that given a test input *x* chooses a classifier *c* ~ *Q* (the posterior) and returns *c*(*x*)
- We are interested in the relation between two quantities:

$$Q_{\mathcal{D}} = \mathbb{E}_{\boldsymbol{c} \sim \boldsymbol{Q}}[\boldsymbol{c}_{\mathcal{D}}]$$

the true error rate of the probabilistic classifier and

$$\hat{\textit{Q}}_{\mathcal{S}} = \mathbb{E}_{\textit{c}\sim\textit{Q}}[\hat{\textit{c}}_{\mathcal{S}}]$$

A (10) + A (10) +

its empirical error rate

Definitions PAC-Bayes Theorem Proof outline Applications

Definitions for main result Generalisation error

Note that this does not bound the posterior average but we have

$$\Pr_{(x,y)\sim\mathcal{D}}(\operatorname{sgn}(\mathbb{E}_{c\sim Q}[c(x)])\neq y)\leq 2Q_{\mathcal{D}}.$$

since for any point *x* misclassified by sgn ($\mathbb{E}_{c\sim Q}[c(x)]$) the probability of a random $c \sim Q$ misclassifying is at least 0.5.

• □ > • □ > • □ > • □ > • □ > ·

Definitions PAC-Bayes Theorem Proof outline Applications

PAC-Bayes Theorem

 Fix an arbitrary D, arbitrary prior P, and confidence δ, then with probability at least 1 – δ over samples S ~ D^m, all posteriors Q satisfy

$$\operatorname{KL}(\hat{Q}_{\mathcal{S}} \| Q_{\mathcal{D}}) \leq \frac{\operatorname{KL}(Q \| P) + \ln((m+1)/\delta)}{m}$$

where KL is the KL divergence between distributions

$$\operatorname{KL}(Q \| P) = \mathbb{E}_{c \sim Q} \left[\ln \frac{Q(c)}{P(c)} \right]$$

(日)

with \hat{Q}_{S} and Q_{D} considered as distributions on $\{0, +1\}$.

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (1/3)

$$\Pr_{\mathcal{S} \sim \mathcal{D}^m} \left(\mathbb{E}_{\boldsymbol{c} \sim \mathcal{P}} \frac{1}{\Pr_{\mathcal{S}' \sim \mathcal{D}^m} (\hat{\boldsymbol{c}}_{\mathcal{S}} = \hat{\boldsymbol{c}}_{\mathcal{S}'})} \leq \frac{m+1}{\delta} \right) \geq 1 - \delta$$

• This follows from considering the expectation divided into probability of particular empirical error for any *c*:

$$\mathbb{E}_{S\sim\mathcal{D}^m}\frac{1}{\Pr_{S'\sim\mathcal{D}^m}(\hat{c}_S=\hat{c}_{S'})}=\sum_k Pr_{S\sim\mathcal{D}^m}(\hat{c}_S=k)\frac{1}{\Pr_{S'\sim\mathcal{D}^m}(\hat{c}_{S'}=k)}=m+1.$$

Taking expectations wrt to c and reversing the expectations

$$\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{c \sim \mathcal{P}} \frac{1}{\Pr_{S' \sim \mathcal{D}^m} (\hat{c}_S = \hat{c}_{S'})} = m + 1$$

and the result follows from Markov's inequality

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (1/3)

$$\Pr_{\mathcal{S}\sim\mathcal{D}^{m}}\left(\mathbb{E}_{\boldsymbol{c}\sim\boldsymbol{P}}\frac{1}{\Pr_{\mathcal{S}'\sim\mathcal{D}^{m}}(\hat{\boldsymbol{c}}_{\mathcal{S}}=\hat{\boldsymbol{c}}_{\mathcal{S}'})}\leq\frac{m+1}{\delta}\right)\geq1-\delta$$

• This follows from considering the expectation divided into probability of particular empirical error for any *c*:

$$\mathbb{E}_{S\sim\mathcal{D}^m}\frac{1}{\Pr_{S'\sim\mathcal{D}^m}(\hat{c}_S=\hat{c}_{S'})}=\sum_k Pr_{S\sim\mathcal{D}^m}(\hat{c}_S=k)\frac{1}{\Pr_{S'\sim\mathcal{D}^m}(\hat{c}_{S'}=k)}=m+1.$$

Taking expectations wrt to c and reversing the expectations

$$\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{c \sim P} \frac{1}{\Pr_{S' \sim \mathcal{D}^m} (\hat{c}_S = \hat{c}_{S'})} = m + 1$$

and the result follows from Markov's inequality.

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (2/3)

1

 $\frac{\mathbb{E}_{c \sim Q} \ln \frac{1}{\Pr_{S' \sim \mathcal{D}^{m}}(\hat{c}_{S} = \hat{c}_{S'})}}{KL(\hat{Q}_{S} \| Q_{\mathcal{D}})} > KL(\hat{Q}_{S} \| Q_{\mathcal{D}})$ m

 This follows by considering the probabilities that the two empirical estimates are equal, applying the relative entropy Chernoff bound and then using the concavity of the KL divergence as a function of both arguments.

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (2/3)

1

$$\frac{\mathbb{E}_{c \sim Q} \ln \frac{1}{\Pr_{S' \sim \mathcal{D}^{m}}(\hat{c}_{S} = \hat{c}_{S'})}}{m} \geq \mathrm{KL}(\hat{Q}_{S} \| Q_{\mathcal{D}})$$

 This follows by considering the probabilities that the two empirical estimates are equal, applying the relative entropy Chernoff bound and then using the concavity of the KL divergence as a function of both arguments.

• □ > • □ > • □ > • □ > • □ > ·

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (3/3)

Consider the distribution

$$egin{aligned} & P_G(c) = rac{1}{ ext{Pr}_{S' \sim \mathcal{D}^m}(\hat{c}_{S'} = \hat{c}_S) \mathbb{E}_{d \sim P} rac{1}{ ext{Pr}_{S' \sim \mathcal{D}^m}(\hat{d}_S = \hat{d}_{S'})}} P(c) \end{aligned}$$

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

<ロ> <同> <同> < 同> < 同> < 同> <

크

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (2/3)

$$\begin{array}{lll} 0 & \leq & \operatorname{KL}(Q \| P_G) \\ & = & \operatorname{KL}(Q \| P) - \mathbb{E}_{c \sim Q} \ln \frac{1}{\operatorname{Pr}_{S' \sim \mathcal{D}^m}(\hat{c}_{S'} = \hat{c}_S)} \\ & & + \ln \mathbb{E}_{d \sim P} \frac{1}{\operatorname{Pr}_{S' \sim \mathcal{D}^m}(\hat{d}_S = \hat{d}_{S'})} \end{array}$$

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

æ

Definitions PAC-Bayes Theorem Proof outline Applications

Ingredients of proof (3/3)

$$\begin{split} m \mathrm{KL}(\hat{Q}_{S} \| Q_{\mathcal{D}}) &\leq \mathbb{E}_{c \sim Q} \ln \frac{1}{\Pr_{S' \sim \mathcal{D}^{m}}(\hat{c}_{S'} = \hat{c}_{S})} \\ &\leq \mathrm{KL}(Q \| P) + \ln \mathbb{E}_{d \sim P} \frac{1}{\Pr_{S' \sim \mathcal{D}^{m}}(\hat{d}_{S} = \hat{d}_{S'})} \\ &\leq \mathrm{KL}(Q \| P) + \frac{m+1}{\delta} \end{split}$$

with probability greater than $1 - \delta$.

・ロト ・聞 ト ・ ヨ ト ・ ヨ ト

크

Definitions PAC-Bayes Theorem Proof outline Applications

Finite Classes

If we take a finite class of functions h₁,..., h_N with prior distribution p₁,..., p_N and assume that the posterior is concentrated on a single function, the generalisation is bounded by

$$\mathrm{KL}(\widehat{\mathrm{err}}(h_i) \| \mathrm{err}(h_i)) \leq \frac{-\log(p_i) + \ln((m+1)/\delta)}{m}$$

 This is the standard result for finite classes with the slight refinement that it involves the KL divergence between empirical and true error and the extra log(m + 1) term on the rhs.

Definitions PAC-Bayes Theorem Proof outline Applications

Finite Classes

If we take a finite class of functions h₁,..., h_N with prior distribution p₁,..., p_N and assume that the posterior is concentrated on a single function, the generalisation is bounded by

$$\operatorname{KL}(\widehat{\operatorname{err}}(h_i) \| \operatorname{err}(h_i)) \leq rac{-\log(p_i) + \ln((m+1)/\delta)}{m}$$

 This is the standard result for finite classes with the slight refinement that it involves the KL divergence between empirical and true error and the extra log(m + 1) term on the rhs.

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Definitions PAC-Bayes Theorem Proof outline Applications

Other extensions/applications

- Matthias Seeger developed the theory for bounding the error of a Gaussian process classifier.
- Olivier Catoni has extended the result to exchangeable distributions enabling him to get a PAC-Bayes version of Vapnik-Chervonenkis bounds.
- Germain et al have extended to more general loss functions than just binary.
- David McAllester has extended the approach to structured output learning.

Definitions PAC-Bayes Theorem Proof outline Applications

Other extensions/applications

- Matthias Seeger developed the theory for bounding the error of a Gaussian process classifier.
- Olivier Catoni has extended the result to exchangeable distributions enabling him to get a PAC-Bayes version of Vapnik-Chervonenkis bounds.
- Germain et al have extended to more general loss functions than just binary.
- David McAllester has extended the approach to structured output learning.

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Definitions PAC-Bayes Theorem Proof outline Applications

Other extensions/applications

- Matthias Seeger developed the theory for bounding the error of a Gaussian process classifier.
- Olivier Catoni has extended the result to exchangeable distributions enabling him to get a PAC-Bayes version of Vapnik-Chervonenkis bounds.
- Germain et al have extended to more general loss functions than just binary.
- David McAllester has extended the approach to structured output learning.

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Definitions PAC-Bayes Theorem Proof outline Applications

Other extensions/applications

- Matthias Seeger developed the theory for bounding the error of a Gaussian process classifier.
- Olivier Catoni has extended the result to exchangeable distributions enabling him to get a PAC-Bayes version of Vapnik-Chervonenkis bounds.
- Germain et al have extended to more general loss functions than just binary.
- David McAllester has extended the approach to structured output learning.

A (B) > A (B) > A (B) >

Definitions PAC-Bayes Theorem Proof outline Applications

Linear classifiers and SVMs

Focus in on linear function application (Langford & ST)

- How the application is made
- Extensions to learning the prior
- Some results on UCI datasets to give an idea of what can be achieved

Definitions PAC-Bayes Theorem Proof outline Applications

Linear classifiers and SVMs

- Focus in on linear function application (Langford & ST)
- How the application is made
- Extensions to learning the prior
- Some results on UCI datasets to give an idea of what can be achieved
Definitions PAC-Bayes Theorem Proof outline Applications

Linear classifiers and SVMs

- Focus in on linear function application (Langford & ST)
- How the application is made
- Extensions to learning the prior
- Some results on UCI datasets to give an idea of what can be achieved

(日)

Definitions PAC-Bayes Theorem Proof outline Applications

Linear classifiers and SVMs

- Focus in on linear function application (Langford & ST)
- How the application is made
- Extensions to learning the prior
- Some results on UCI datasets to give an idea of what can be achieved

(日)

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Linear classifiers

- We will choose the prior and posterior distributions to be Gaussians with unit variance.
- The prior P will be centered at the origin with unit variance
- The specification of the centre for the posterior Q(w, μ) will be by a unit vector w and a scale factor μ.

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 一日 ト ・ 日 ト

Linear classifiers

- We will choose the prior and posterior distributions to be Gaussians with unit variance.
- The prior P will be centered at the origin with unit variance
- The specification of the centre for the posterior Q(w, μ) will be by a unit vector w and a scale factor μ.

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Linear classifiers

- We will choose the prior and posterior distributions to be Gaussians with unit variance.
- The prior P will be centered at the origin with unit variance
- The specification of the centre for the posterior Q(w, μ) will be by a unit vector w and a scale factor μ.

General Approach Learning the prior New prior for linear functions Prior-SVM

< 17 ▶

PAC-Bayes Bound for SVM (1/2)



General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (1/2)



• **Prior** *P* is Gaussian $\mathcal{N}(0, 1)$

< 17 ▶

Posterior is in the direction w

- ۲
 - 1

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (1/2)



- **Prior** *P* is Gaussian $\mathcal{N}(0, 1)$
- Posterior is in the direction w
- at **distance** μ from the origin

< 17 >

()

۲

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (1/2)



- **Prior** *P* is Gaussian $\mathcal{N}(0, 1)$
- Posterior is in the direction w
- at **distance** μ from the origin
- **Posterior** *Q* is Gaussian

< 17 ▶

A B F A B F

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- $Q_D(\mathbf{w}, \mu)$ true performance of the stochastic classifier
- SVM is deterministic classifier that exactly corresponds to sgn (E_{c~Q(w,μ)}[c(x)]) as centre of the Gaussian gives the same classification as halfspace with more weight.
- Hence its error bounded by 2Q_D(**w**, μ), since as observed above if x misclassified at least half of c ~ Q err.

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- $Q_D(\mathbf{w}, \mu)$ true performance of the stochastic classifier
- SVM is deterministic classifier that exactly corresponds to sgn (E_{c~Q(w,μ)}[c(x)]) as centre of the Gaussian gives the same classification as halfspace with more weight.
- Hence its error bounded by 2Q_D(**w**, μ), since as observed above if x misclassified at least half of c ~ Q err.

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- $Q_D(\mathbf{w}, \mu)$ true performance of the stochastic classifier
- SVM is deterministic classifier that exactly corresponds to sgn (E_{c~Q(w,μ)}[c(x)]) as centre of the Gaussian gives the same classification as halfspace with more weight.
- Hence its error bounded by 2Q_D(**w**, μ), since as observed above if x misclassified at least half of c ~ Q err.

General Approach Learning the prior New prior for linear functions Prior-SVM

• □ ▶ • □ ▶ • □ ▶ • □ ▶ •

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- $Q_D(\mathbf{w}, \mu)$ true performance of the stochastic classifier
- SVM is deterministic classifier that exactly corresponds to sgn (E_{c~Q(w,μ)}[c(x)]) as centre of the Gaussian gives the same classification as halfspace with more weight.
- Hence its error bounded by 2Q_D(**w**, μ), since as observed above if x misclassified at least half of c ~ Q err.

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

PAC-Bayes Bound for SVM (2/2)

Linear classifiers performance may be bounded by

$$\mathsf{KL}(\left[\hat{\boldsymbol{Q}}_{\mathcal{S}}(\boldsymbol{w},\mu)\right] \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

• $\hat{Q}_{S}(\mathbf{w},\mu)$ stochastic measure of the training error

$$\hat{Q}_{S}(\boldsymbol{w},\mu) = \mathbb{E}_{m}[\tilde{F}(\mu\gamma(\boldsymbol{x},y))]$$
$$\gamma(\boldsymbol{x},y) = (\boldsymbol{y}\boldsymbol{w}^{T}\phi(\boldsymbol{x}))/(\|\phi(\boldsymbol{x})\|\|\boldsymbol{w}\|)$$
$$\tilde{F}(t) = 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t} e^{-x^{2}/2} dx$$

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (2/2)

Linear classifiers performance may be bounded by

$$\mathsf{KL}(\left[\hat{\boldsymbol{Q}}_{\mathcal{S}}(\boldsymbol{w},\mu)\right] \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(\boldsymbol{P} \| \boldsymbol{Q}(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

Q̂_S(w, μ) stochastic measure of the training error

$$\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) = \mathbb{E}_{m}[\tilde{F}(\mu\gamma(\boldsymbol{x},\boldsymbol{y}))]$$
$$\gamma(\boldsymbol{x},\boldsymbol{y}) = (\boldsymbol{y}\boldsymbol{w}^{T}\phi(\boldsymbol{x}))/(\|\phi(\boldsymbol{x})\|\|\boldsymbol{w}\|)$$
$$\tilde{F}(t) = 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t} e^{-x^{2}/2} dx$$

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(P \| Q(\boldsymbol{w},\mu))}{m} + \ln \frac{m+1}{\delta}$$

- Prior $P \equiv$ Gaussian centered on the origin
- Posterior $Q \equiv$ Gaussian along **w** at a distance μ from the origin

•
$$KL(P||Q) = \mu^2/2$$

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (2/2)

Linear classifiers performance may be bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(P \| Q(\boldsymbol{w},\mu))}{m} + \ln \frac{m+1}{\delta}$$

• Prior $P \equiv$ Gaussian centered on the origin

• Posterior $Q \equiv$ Gaussian along **w** at a distance μ from the origin

•
$$KL(P||Q) = \mu^2/2$$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{|\mathsf{KL}(\boldsymbol{P}\| Q(\boldsymbol{w},\mu))| + \ln \frac{m+1}{\delta}}{m}$$

- Prior $P \equiv$ Gaussian centered on the origin
- Posterior $Q \equiv$ Gaussian along **w** at a distance μ from the origin

•
$$KL(P||Q) = \mu^2/2$$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \mathcal{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{|\mathsf{KL}(\boldsymbol{P}\| \mathcal{Q}(\boldsymbol{w},\mu))| + \ln \frac{m+1}{\delta}}{m}$$

- Prior $P \equiv$ Gaussian centered on the origin
- Posterior $Q \equiv$ Gaussian along **w** at a distance μ from the origin

•
$$KL(P||Q) = \mu^2/2$$

General Approach Learning the prior New prior for linear functions Prior-SVM

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(P \| Q(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- δ is the confidence
- The bound holds with probability 1δ over the random i.i.d. selection of the training data.

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

Linear classifiers performance may be bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(P \| Q(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

• δ is the confidence

• The bound holds with probability $1 - \delta$ over the random i.i.d. selection of the training data.

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

PAC-Bayes Bound for SVM (2/2)

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{\mathsf{KL}(P \| Q(\boldsymbol{w},\mu)) + \ln \frac{m+1}{\delta}}{m}$$

- δ is the confidence
- The bound holds with probability 1δ over the random i.i.d. selection of the training data.

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Learning the prior (1/3)

Bound depends on the distance between prior and posterior

- Better prior (closer to posterior) would lead to tighter bound
- Learn the prior P with part of the data
- Introduce the learnt prior in the bound
- Compute stochastic error with remaining data

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

- Bound depends on the distance between prior and posterior
- Better prior (closer to posterior) would lead to tighter bound
- Learn the prior P with part of the data
- Introduce the learnt prior in the bound
- Compute stochastic error with remaining data

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

- Bound depends on the distance between prior and posterior
- Better prior (closer to posterior) would lead to tighter bound
- Learn the prior P with part of the data
- Introduce the learnt prior in the bound
- Compute stochastic error with remaining data

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

- Bound depends on the distance between prior and posterior
- Better prior (closer to posterior) would lead to tighter bound
- Learn the prior P with part of the data
- Introduce the learnt prior in the bound
- Compute stochastic error with remaining data

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

- Bound depends on the distance between prior and posterior
- Better prior (closer to posterior) would lead to tighter bound
- Learn the prior P with part of the data
- Introduce the learnt prior in the bound
- Compute stochastic error with remaining data

General Approach Learning the prior New prior for linear functions Prior-SVM

New prior for the SVM (3/3)



John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

General Approach Learning the prior New prior for linear functions Prior-SVM

New prior for the SVM (3/3)



Solve SVM with subset of patterns

• Prior in the **direction w**_r

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

General Approach Learning the prior New prior for linear functions Prior-SVM

New prior for the SVM (3/3)



- Solve SVM with subset of patterns
- Prior in the direction w_r
- Posterior like PAC-Bayes Bound

(日)

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

General Approach Learning the prior New prior for linear functions Prior-SVM

New prior for the SVM (3/3)



- Solve SVM with subset of patterns
- Prior in the **direction w**_r
- Posterior like PAC-Bayes Bound
- New bound proportional to KL(P||Q)

(日)

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 四 ト ・ 回 ト ・ 回 ト

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2} + \ln \frac{(m-r+1)J}{\delta}}{m-r}$$

• $Q_{\mathcal{D}}(\boldsymbol{w},\mu)$ true performance of the classifier

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 四 ト ・ 回 ト ・ 回 ト

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| \mathbf{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \|\mu \boldsymbol{w} - \eta \boldsymbol{w}_{r}\|^{2} + \ln \frac{(m-r+1)J}{\delta}}{m-r}$$

• $Q_{\mathcal{D}}(\mathbf{w},\mu)$ true performance of the classifier

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\widehat{\boldsymbol{Q}_{\mathcal{S}}(\boldsymbol{w},\mu)} \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2} + \ln \frac{(m-r+1)J}{\delta}}{m-r}$$

• $\hat{Q}_{S}(w, \mu)$ stochastic measure of the training error on remaining data

$$\hat{Q}(\boldsymbol{w},\mu)_{S} = \mathbb{E}_{\boldsymbol{m}-\boldsymbol{r}}[\tilde{F}(\mu\gamma(\boldsymbol{x},\boldsymbol{y}))]$$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロ・ ・ 四・ ・ 回・ ・ 回・

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\left[\hat{\boldsymbol{Q}}_{\mathcal{S}}(\boldsymbol{w},\mu)\right] \| \boldsymbol{Q}_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2} + \ln \frac{(m-r+1)J}{\delta}}{m-r}$$

• $\hat{Q}_{S}(w, \mu)$ stochastic measure of the training error on remaining data

$$\hat{Q}(\boldsymbol{w},\mu)_{\mathcal{S}} = \mathbb{E}_{\boldsymbol{m}-\boldsymbol{r}}[\tilde{F}(\mu\gamma(\boldsymbol{x},\boldsymbol{y}))]$$

General Approach Learning the prior New prior for linear functions Prior-SVM

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

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2}}{m-r} + \ln \frac{(m-r+1)J}{\delta}$$

• $0.5 \|\mu \mathbf{w} - \eta \mathbf{w}_r\|^2$ distance between prior and posterior
General Approach Learning the prior New prior for linear functions Prior-SVM

<ロ> <同> <同> < 同> < 同> < 同> < □> <

э

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2}}{m-r} + \ln \frac{(m-r+1)J}{\delta}$$

• $0.5 \|\mu \mathbf{w} - \eta \mathbf{w}_r\|^2$ distance between prior and posterior

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 四 ト ・ 回 ト ・ 回 ト

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_{r} \|^{2} + \ln \frac{(m-r+1)J}{\delta}}{[m-r]}$$

• Penalty term only dependent on the remaining data m - r

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・聞 ト ・ ヨ ト ・ ヨ ト

New Bound for the SVM (2/3)

SVM performance may be tightly bounded by

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}}(\boldsymbol{w},\mu) \| Q_{\mathcal{D}}(\boldsymbol{w},\mu)) \leq \frac{0.5 \| \mu \boldsymbol{w} - \eta \boldsymbol{w}_r \|^2 + \ln \frac{(m-r+1)J}{\delta}}{\left[\frac{m-r}{\delta} \right]}$$

Penalty term only dependent on the remaining data m - r

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 四 ト ・ 回 ト ・ 回 ト

Prior-SVM

- New bound proportional to $\|\mu \mathbf{w} \eta \mathbf{w}_r\|^2$
- Classifier that optimises the bound
- Optimisation problem to determine the p-SVM

$$\min_{\boldsymbol{W},\xi_i} \left[\frac{1}{2} \| \boldsymbol{w} - \boldsymbol{w}_r \|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

s.t. $y_i \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) \ge 1 - \xi_i$ $i = 1, \dots, m-r$
 $\xi_i \ge 0$ $i = 1, \dots, m-r$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

Prior-SVM

- New bound proportional to $\|\mu \mathbf{w} \eta \mathbf{w}_r\|^2$
- Classifier that optimises the bound
- Optimisation problem to determine the p-SVM

$$\min_{\boldsymbol{W},\xi_i} \left[\frac{1}{2} \| \boldsymbol{w} - \boldsymbol{w}_r \|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

s.t. $y_i \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_i) \ge 1 - \xi_i$ $i = 1, \dots, m-r$
 $\xi_i \ge 0$ $i = 1, \dots, m-r$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Prior-SVM

- New bound proportional to $\|\mu \mathbf{w} \eta \mathbf{w}_r\|^2$
- Classifier that optimises the bound
- Optimisation problem to determine the p-SVM

$$\min_{\boldsymbol{W},\xi_i} \begin{bmatrix} \frac{1}{2} \| \boldsymbol{w} - \boldsymbol{w}_r \|^2 + C \sum_{i=1}^{m-r} \xi_i \end{bmatrix}$$

s.t. $y_i \boldsymbol{w}^T \phi(\boldsymbol{x}_i) \ge 1 - \xi_i$ $i = 1, \dots, m-r$
 $\xi_i \ge 0$ $i = 1, \dots, m-r$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

Prior-SVM

- New bound proportional to $\|\mu \mathbf{w} \eta \mathbf{w}_r\|^2$
- Classifier that optimises the bound
- Optimisation problem to determine the p-SVM

$$\min_{\boldsymbol{W},\xi_i} \left[\frac{1}{2} \| \boldsymbol{w} - \boldsymbol{w}_r \|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

s.t. $y_i \boldsymbol{w}^T \phi(\boldsymbol{x}_i) \ge 1 - \xi_i$ $i = 1, \dots, m-r$
 $\xi_i \ge 0$ $i = 1, \dots, m-r$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 一日 ト ・ 日 ト

Bound for p-SVM

- Determine the prior with a subset of the training examples to obtain w_r
- Solve p-SVM and obtain w
- Imagin for the stochastic classifier \hat{Q}_s

$$\gamma(\mathbf{x}_j, y_j) = \frac{y_j \mathbf{w}^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_j)\| \|\mathbf{w}\|} \qquad j = 1, \dots, m - r$$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 一日 ト ・ 日 ト

Bound for p-SVM

- Determine the prior with a subset of the training examples to obtain w_r
- Solve p-SVM and obtain w
- 3 Margin for the stochastic classifier \hat{Q}_s

$$\gamma(\mathbf{x}_j, y_j) = \frac{y_j \mathbf{w}^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_j)\| \|\mathbf{w}\|} \qquad j = 1, \dots, m - r$$

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Bound for p-SVM

- Determine the prior with a subset of the training examples to obtain w_r
- Solve p-SVM and obtain w
- **O Margin** for the stochastic classifier \hat{Q}_s

$$\gamma(\boldsymbol{x}_j, y_j) = \frac{y_j \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_j)}{\|\boldsymbol{\phi}(\boldsymbol{x}_j)\| \|\boldsymbol{w}\|} \qquad j = 1, \dots, m - r$$

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Bound for p-SVM

- Determine the prior with a subset of the training examples to obtain w_r
- Solve p-SVM and obtain w
- **O Margin** for the stochastic classifier \hat{Q}_s

$$\gamma(\boldsymbol{x}_j, \boldsymbol{y}_j) = \frac{y_j \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_j)}{\|\boldsymbol{\phi}(\boldsymbol{x}_j)\| \|\boldsymbol{w}\|} \qquad j = 1, \dots, m - r$$

General Approach Learning the prior New prior for linear functions Prior-SVM

η -Prior-SVM

- Consider using a prior distribution *P* that is elongated in the direction of w_r
- This will mean that there is low penalty for large projections onto this direction
- Translates into an optimisation:

$$\min_{\boldsymbol{v},\eta,\xi_i} \left[\frac{1}{2} \|\boldsymbol{v}\|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

subject to

$$y_i(\mathbf{v} + \eta \mathbf{w}_r)^T \phi(\mathbf{x}_i) \ge 1 - \xi_i \qquad i = 1, \dots, m - r$$

$$\xi_i \ge 0 \qquad i = 1, \dots, m - r$$

・ロト ・ 四 ト ・ 回 ト ・ 回 ト

General Approach Learning the prior New prior for linear functions Prior-SVM

η -Prior-SVM

- Consider using a prior distribution *P* that is elongated in the direction of w_r
- This will mean that there is low penalty for large projections onto this direction
- Translates into an optimisation:

$$\min_{\boldsymbol{\nu},\eta,\xi_i}\left[\frac{1}{2}\|\boldsymbol{\nu}\|^2+C\sum_{i=1}^{m-r}\xi_i\right]$$

subject to

$$y_i(\mathbf{v} + \eta \mathbf{w}_r)^T \phi(\mathbf{x}_i) \ge 1 - \xi_i \qquad i = 1, \dots, m - r$$

$$\xi_i \ge 0 \qquad i = 1, \dots, m - r$$

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

 Outline
 General Approach

 Links
 Learning the prior

 PAC-Bayes Analysis
 New prior for linear functions

 Linear Classifiers
 Prior-SVM

η -Prior-SVM

- Consider using a prior distribution *P* that is elongated in the direction of w_r
- This will mean that there is low penalty for large projections onto this direction
- Translates into an optimisation:

$$\min_{\mathbf{v},\eta,\xi_i} \left[\frac{1}{2} \|\mathbf{v}\|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

subject to

$$y_i(\mathbf{v} + \eta \mathbf{w}_r)^T \phi(\mathbf{x}_i) \ge 1 - \xi_i \qquad i = 1, \dots, m - r$$

$$\xi_i \ge 0 \qquad i = 1, \dots, m - r$$

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

 Outline
 General Approach

 Links
 Learning the prior

 PAC-Bayes Analysis
 New prior for linear functions

 Linear Classifiers
 Prior-SVM

η -Prior-SVM

- Consider using a prior distribution *P* that is elongated in the direction of w_r
- This will mean that there is low penalty for large projections onto this direction
- Translates into an optimisation:

$$\min_{\mathbf{v},\eta,\xi_i} \left[\frac{1}{2} \|\mathbf{v}\|^2 + C \sum_{i=1}^{m-r} \xi_i \right]$$

subject to

$$y_i(\mathbf{v} + \eta \mathbf{w}_r)^T \phi(\mathbf{x}_i) \ge 1 - \xi_i \qquad i = 1, \dots, m - r$$

$$\xi_i \ge 0 \qquad i = 1, \dots, m - r$$

(日)

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Bound for η -prior-SVM

- Prior is elongated along the line of w_r but spherical with variance 1 in other directions
- Posterior again on the line of w at a distance μ chosen to optimise the bound.
- Resulting bound depends on a benign parameter *τ* determining the variance in the direction **w**_r

$$\mathsf{KL}(\hat{Q}_{S\setminus R}(\mathbf{w},\mu) \| \mathcal{Q}_{\mathcal{D}}(\mathbf{w},\mu)) \leq \frac{0.5(\ln(\tau^2) + \tau^{-2} - 1 + P_{\mathbf{w}_r}^{\parallel}(\mu\mathbf{w} - \mathbf{w}_r)^2/\tau^2 + P_{\mathbf{w}_r}^{\perp}(\mu\mathbf{w})^2) + \ln(\frac{m-r+1}{\delta})}{m-r}$$

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Bound for η -prior-SVM

- Prior is elongated along the line of w_r but spherical with variance 1 in other directions
- Posterior again on the line of w at a distance μ chosen to optimise the bound.
- Resulting bound depends on a benign parameter *τ* determining the variance in the direction **w**_r

$$\mathsf{KL}(\hat{Q}_{S\setminus R}(\mathbf{w},\mu) \| \mathcal{Q}_{\mathcal{D}}(\mathbf{w},\mu)) \leq \frac{0.5(\ln(\tau^2) + \tau^{-2} - 1 + P_{\mathbf{w}_r}^{\parallel}(\mu\mathbf{w} - \mathbf{w}_r)^2/\tau^2 + P_{\mathbf{w}_r}^{\perp}(\mu\mathbf{w})^2) + \ln(\frac{m-r+1}{\delta})}{m-r}$$

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Bound for η -prior-SVM

- Prior is elongated along the line of w_r but spherical with variance 1 in other directions
- Posterior again on the line of w at a distance μ chosen to optimise the bound.
- Resulting bound depends on a benign parameter *τ* determining the variance in the direction **w**_r

$$\mathsf{KL}(\hat{Q}_{\mathcal{S}\setminus R}(\mathbf{w},\mu) \| \mathcal{Q}_{\mathcal{D}}(\mathbf{w},\mu)) \leq \\ \frac{0.5(\ln(\tau^2) + \tau^{-2} - 1 + \mathcal{P}_{\mathbf{w}_r}^{\parallel}(\mu\mathbf{w} - \mathbf{w}_r)^2/\tau^2 + \mathcal{P}_{\mathbf{w}_r}^{\perp}(\mu\mathbf{w})^2) + \ln(\frac{m-r+1}{\delta})}{m-r}$$

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Comparison with X-fold Xvalidation, PAC-Bayes Bound and the Prior PAC-Bayes Bound
- UCI datasets
- Select C and σ that lead to minimum Classification Error (CE)
 - For X-F XV select the pair that minimize the validation error
 - For PAC-Bayes Bound and Prior PAC-Bayes Bound select the pair that minimize the bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Comparison with X-fold Xvalidation, PAC-Bayes Bound and the Prior PAC-Bayes Bound
- UCI datasets
- Select C and σ that lead to minimum Classification Error (CE)
 - For X-F XV select the pair that minimize the validation error
 - For PAC-Bayes Bound and Prior PAC-Bayes Bound select the pair that minimize the bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Comparison with X-fold Xvalidation, PAC-Bayes Bound and the Prior PAC-Bayes Bound
- UCI datasets
- Select C and σ that lead to minimum Classification Error (CE)
 - For X-F XV select the pair that minimize the validation error
 - For PAC-Bayes Bound and Prior PAC-Bayes Bound select the pair that minimize the bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Comparison with X-fold Xvalidation, PAC-Bayes Bound and the Prior PAC-Bayes Bound
- UCI datasets
- Select C and σ that lead to minimum Classification Error (CE)
 - For X-F XV select the pair that minimize the validation error
 - For PAC-Bayes Bound and Prior PAC-Bayes Bound select the pair that minimize the bound

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

- Comparison with X-fold Xvalidation, PAC-Bayes Bound and the Prior PAC-Bayes Bound
- UCI datasets
- Select C and σ that lead to minimum Classification Error (CE)
 - For X-F XV select the pair that minimize the validation error
 - For PAC-Bayes Bound and Prior PAC-Bayes Bound select the pair that minimize the bound

General Approach Learning the prior New prior for linear functions Prior-SVM

(日)

Description of the Datasets

Problem	# samples	input dim.	Pos/Neg	
Handwritten-digits	5620	64	2791 / 2829	
Waveform	5000	21	1647 / 3353	
Pima	768	8	268 / 500	
Ringnorm	7400	20	3664 / 3736	
Spam	4601	57	1813 / 2788	

Table: Description of datasets in terms of number of patterns, number of input variables and number of positive/negative examples.

General Approach Learning the prior New prior for linear functions Prior-SVM

3

Results

		Classifier					
		SVM				η Prior SVM	
Problem		2FCV	10FCV	PAC	PrPAC	PrPAC	τ -PrPAC
digits	Bound	-	-	0.175	0.107	0.050	0.047
	CE	0.007	0.007	0.007	0.014	0.010	0.009
waveform	Bound	_	-	0.203	0.185	0.178	0.176
	CE	0.090	0.086	0.084	0.088	0.087	0.086
pima	Bound	_	-	0.424	0.420	0.428	0.416
	CE	0.244	0.245	0.229	0.229	0.233	0.233
ringnorm	Bound	_	-	0.203	0.110	0.053	0.050
	CE	0.016	0.016	0.018	0.018	0.016	0.016
spam	Bound	-	-	0.254	0.198	0.186	0.178
	CE	0.066	0.063	0.067	0.077	0.070	0.072

John Shawe-Taylor University College London Data Dependent Priors in PAC-Bayes Bounds

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound

General Approach Learning the prior New prior for linear functions Prior-SVM

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Frequentist (PAC) and Bayesian approaches to analysing learning lead to introduction of the PAC-Bayes bound
- Detailed look at the ingredients of the theory
- Application to bound the performance of an SVM
- Investigation of learning of the prior of the distribution of classifiers
- Experiments show the new bound can be tighter ...
- ...And reliable for low cost model selection
- p-SVM and η-p-SVM: classifiers that optimise the new bound