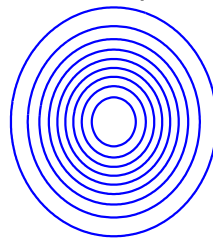


Bayesian Inference: Principles and Practice

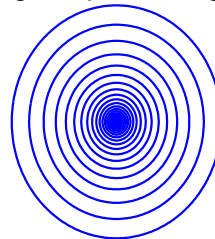
3. Sparse Bayesian Models and the “Relevance Vector Machine”

Mike Tipping

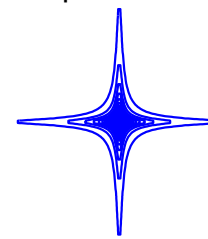
Gaussian prior



Marginal prior: single α



Independent α



Lecture 3: Overview*

- A hierarchical prior for ‘sparse’ linear models
- Sparse regression: (approximate) inference procedure
- Insight into sparsity
- Sparse classification: a further approximation
- Contrast with the support vector machine

*We won't cover all the material in the notes in this lecture ...

Bayes and Contemporary Machine Learning

- We've seen that marginalisation is a valuable component of the Bayesian data modelling paradigm
- We also saw that full Bayesian inference is often analytically intractable, although approximations for simple linear models could be very effective
- Historically, interest in Bayesian “machine learning” (but not statistics!) has focussed on approximations for *non-linear* models, e.g. neural networks:
 - The “evidence procedure” [MacKay]
 - Hybrid Monte Carlo sampling [Neal]
- Good news! — flexible *linear* kernel methods have become very fashionable, thanks mainly to the “support vector machine”

Linear Models and Sparsity

- Much interest in linear models has focused on *sparse* learning algorithms, which set many weights w_m to zero in the function $y(x) = \sum_m w_m \phi_m(x)$
- Sparsity offers:
 - Elegant complexity control
 - Feature extraction
 - Elucidation
- Speed and compactness...



How To Encode Sparsity?

- In the algorithm:

- Heuristic, ‘greedy’, sequential techniques (e.g. “matching pursuit”)

- In the objective function:

$$\hat{E}(\mathbf{w}) = E_{\mathcal{D}}(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

- usually via the penalty term $E_W(\mathbf{w})$

- but can be via the data error term $E_{\mathcal{D}}(\mathbf{w})$ (e.g. the SVM)

- In the prior:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w}) \textcolor{red}{p}(\mathbf{w})}{p(\mathcal{D})}$$

- Recall: correspondence between $E_W(\mathbf{w})$ and $-\log p(\mathbf{w})$

Penalty Terms & Priors

- Most common regularisation term: $E_W(\mathbf{w}) = \sum_{m=1}^M |w_m|^2$
 - Corresponds to a Gaussian prior $p(\mathbf{w}) \propto \exp\left(-\sum_m |w_m|^2\right)$
 - Easy to work with & an effective regulariser, but no sparsity pressure
- ‘Correct’ term would be $E_W(\mathbf{w}) = \sum_m |w_m|^0$, but this is very ugly to work with!
- Instead, $E_W(\mathbf{w}) = \sum_m |w_m|^1$ is a workable compromise which gives reasonable sparsity and reasonable tractability:
 - “LP-machines”
 - “Basis pursuit”
 - “Bayesian Pruning”, as a Laplace prior $p(\mathbf{w}) \propto \exp\left(-\sum_m |w_m|\right)$

A Sparse Bayesian Prior

- In fact, we *can* obtain sparsity by retaining the traditional Gaussian prior
 - This is great news for tractability
 - Furthermore, this approach leads to *scale invariance*
- The modification to our earlier Gaussian prior is subtle:

$$p(\mathbf{w}|\alpha_1, \dots, \alpha_M) = (2\pi)^{-M/2} \prod_{m=1}^M \alpha_m^{1/2} \exp \left\{ -\frac{1}{2} \sum_{m=1}^M \alpha_m w_m^2 \right\}$$

- What's new? We now have M hyperparameters $\alpha = (\alpha_1, \dots, \alpha_M)$, one α_m independently controlling the (inverse) variance of each weight w_m

Hierarchical Priors

- The prior $p(\mathbf{w}|\alpha)$ is Gaussian, but conditioned on α
- We should now define hyperpriors over all α_m
- Previously, we utilised a log-uniform hyperprior — this is a special case of a more general *Gamma* hyperprior
- This gives a *hierarchical* prior: we must marginalise to see its ‘true’ form

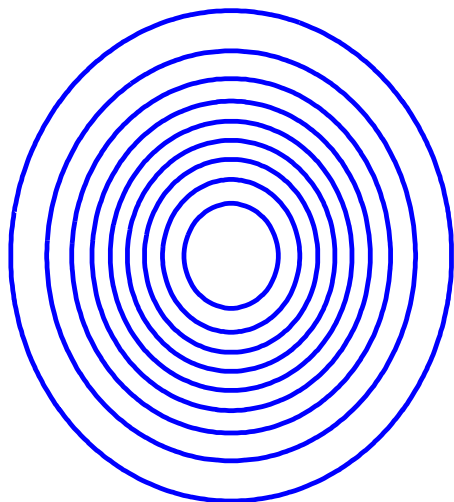
$$p(w_m) = \int p(w_m|\alpha_m) p(\alpha_m) d\alpha_m$$

- For a Gamma $p(\alpha_m)$, $p(w_m)$ is a *Student-t* distribution
- Its equivalent as a penalty function would be: $\sum_m \log |w_m|$

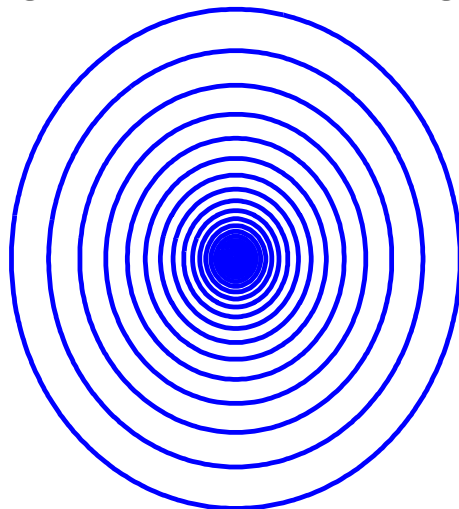
Priors Compared

- Specifying independent hyperparameters α_m is the key to sparsity
- Example marginal priors $p(w_1, w_2)$ illustrated below:

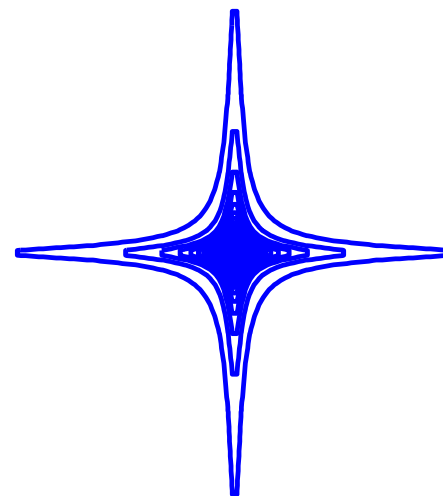
Gaussian prior



Marginal prior: single α



Independent α



A Sparse Bayesian Model for Regression

- As before: independent Gaussian noise: $t_n \sim N(y(\mathbf{x}_n; \mathbf{w}), \sigma^2)$

- This gives a corresponding likelihood:

$$p(\mathbf{t}|\mathbf{w}, \sigma^2) = (2\pi\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{t} - \Phi\mathbf{w}\|^2 \right\}$$

- $\mathbf{t} = (t_1 \dots t_N)^\top$

- $\mathbf{w} = (w_1 \dots w_M)^\top$

- Φ is the $N \times M$ ‘design’ matrix with $\Phi_{nm} = \phi_m(\mathbf{x}_n)$

Approximate Inference

- Desire posterior over all unknowns:

$$p(\mathbf{w}, \alpha, \sigma^2 | \mathbf{t}) = \frac{p(\mathbf{t} | \mathbf{w}, \alpha, \sigma^2) p(\mathbf{w}, \alpha, \sigma^2)}{p(\mathbf{t})}$$

- Can't compute analytically! So as previously, decompose as:

$$p(\mathbf{w}, \alpha, \sigma^2 | \mathbf{t}) \equiv p(\mathbf{w} | \mathbf{t}, \alpha, \sigma^2) p(\alpha, \sigma^2 | \mathbf{t})$$

- $p(\mathbf{w} | \mathbf{t}, \alpha, \sigma^2)$, the 'weight posterior' distribution, is tractable
- $p(\alpha, \sigma^2 | \mathbf{t})$ must be approximated

The Weight Posterior Term

- Given the data, the posterior distribution over weights is Gaussian:

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \alpha, \sigma^2) &= \frac{p(\mathbf{t}|\mathbf{w}, \sigma^2) p(\mathbf{w}|\alpha)}{p(\mathbf{t}|\alpha, \sigma^2)} \\ &= (2\pi)^{-(N+1)/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mu)^\top \Sigma^{-1} (\mathbf{w} - \mu) \right\} \end{aligned}$$

- with

- $\Sigma = (\sigma^{-2} \Phi^\top \Phi + \mathbf{A})^{-1}$

- $\mu = \sigma^{-2} \Sigma \Phi^\top \mathbf{t}$

- $\mathbf{A} = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_M)$

- **Key point:** if $\alpha_m = \infty$, the corresponding $\mu_m = 0$

The Hyperparameter Posterior Term

- As before, for uniform hyperpriors over $\log \alpha$ and $\log \sigma$, $p(\alpha, \sigma^2 | \mathbf{t}) \propto p(\mathbf{t} | \alpha, \sigma^2)$

- Integrate out weights to obtain the *marginal likelihood*:

$$\begin{aligned} p(\mathbf{t} | \alpha, \sigma^2) &= \int p(\mathbf{t} | \mathbf{w}, \sigma^2) p(\mathbf{w} | \alpha) d\mathbf{w}, \\ &= (2\pi)^{-N/2} |\sigma^2 \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^\top|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{t}^\top (\sigma^2 \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^\top)^{-1} \mathbf{t} \right\} \end{aligned}$$

- We maximise $p(\mathbf{t} | \alpha, \sigma^2)$ to find α_{MP} and σ_{MP}^2 (“type-II maximum likelihood”)

- In Lecture 2, we found the single α_{MP} empirically

- For multiple hyperparameters, we must optimise $p(\mathbf{t} | \alpha, \sigma^2)$ directly

Hyperparameter Re-estimation

- Differentiating $\log p(\mathbf{t}|\alpha, \sigma^2)$ gives iterative re-estimation formulae:

$$\alpha_i^{\text{new}} = \frac{\gamma_i}{\mu_i^2}$$

$$(\sigma^2)^{\text{new}} = \frac{\|\mathbf{t} - \Phi\boldsymbol{\mu}\|^2}{N - \sum_{i=1}^M \gamma_i}$$

- For convenience we have defined the quantities

$$\gamma_i = 1 - \alpha_i \Sigma_{ii}$$

- $\gamma_i \in [0, 1]$ is a measure of ‘well-determinedness’ of parameter w_i

Summary of Inference Procedure

- ① Initialise $\{\alpha_j\}$ and σ^2 (or fix latter if known)
- ② Compute weight posterior sufficient statistics μ and Σ
- ③ Compute all $\{\gamma_j\}$, then re-estimate $\{\alpha_j\}$ (and σ^2 if desired)
- ④ Repeat from ② until convergence

- ⑤ Make predictions for new data via the predictive distribution:

$$p(t_*|\mathbf{t}) = \int p(t_*|\mathbf{w}, \sigma_{\text{MP}}^2) p(\mathbf{w}|\mathbf{t}, \alpha_{\text{MP}}, \sigma_{\text{MP}}^2) d\mathbf{w}$$

the mean of which is $y(\mathbf{x}_*; \mu)$

- ⑥ Can ‘delete’ basis functions for which optimal $\alpha_j = \infty$, since $\mu_j = 0$

Sparsity: Algorithmic Insight

- ML(\mathbf{w}): conventional maximum likelihood (or least-squares)

- optimises $M + 1$ parameters, \mathbf{w} and σ^2
- Since $\sigma^2 \rightarrow 0$, $p(\mathbf{t}|\mathbf{w})$ is a δ -function located at \mathbf{t}
- No “Ockham penalty” !

- ML(α): maximising the marginal likelihood

- optimises $M + 1$ parameters, α and σ^2
- $p(\mathbf{t}|\alpha, \sigma^2)$ is a zero-mean *Gaussian process* with covariance

$$\mathbf{C} = \sum_{m=1}^M \alpha_m^{-1} \phi_m \phi_m^\top + \sigma^2 \mathbf{I}$$

where $\phi_m = [\phi_m(\mathbf{x}_1), \dots, \phi_m(\mathbf{x}_N)]^\top$

Observation Space \mathbf{R}^N

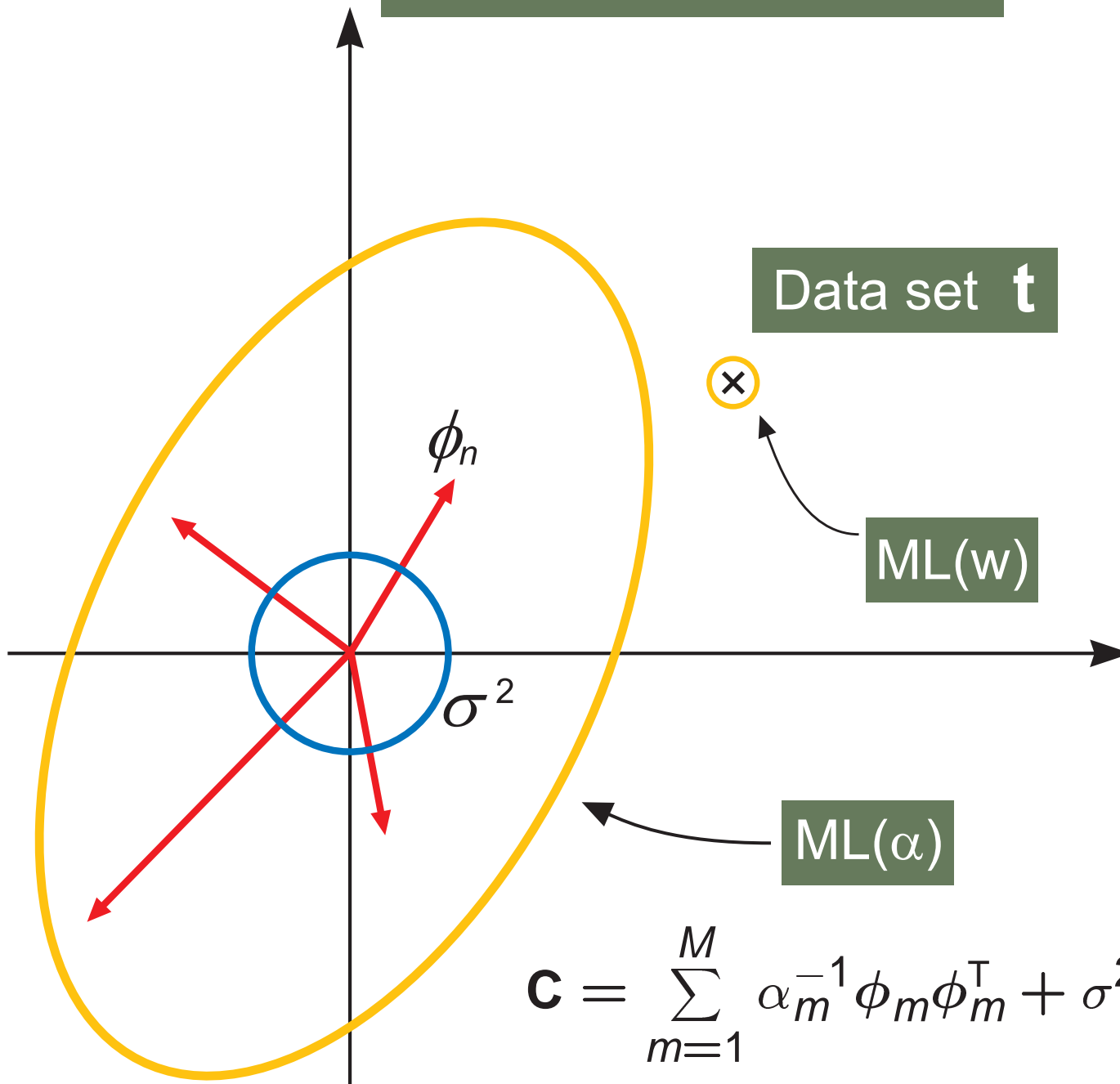
Data set \mathbf{t}



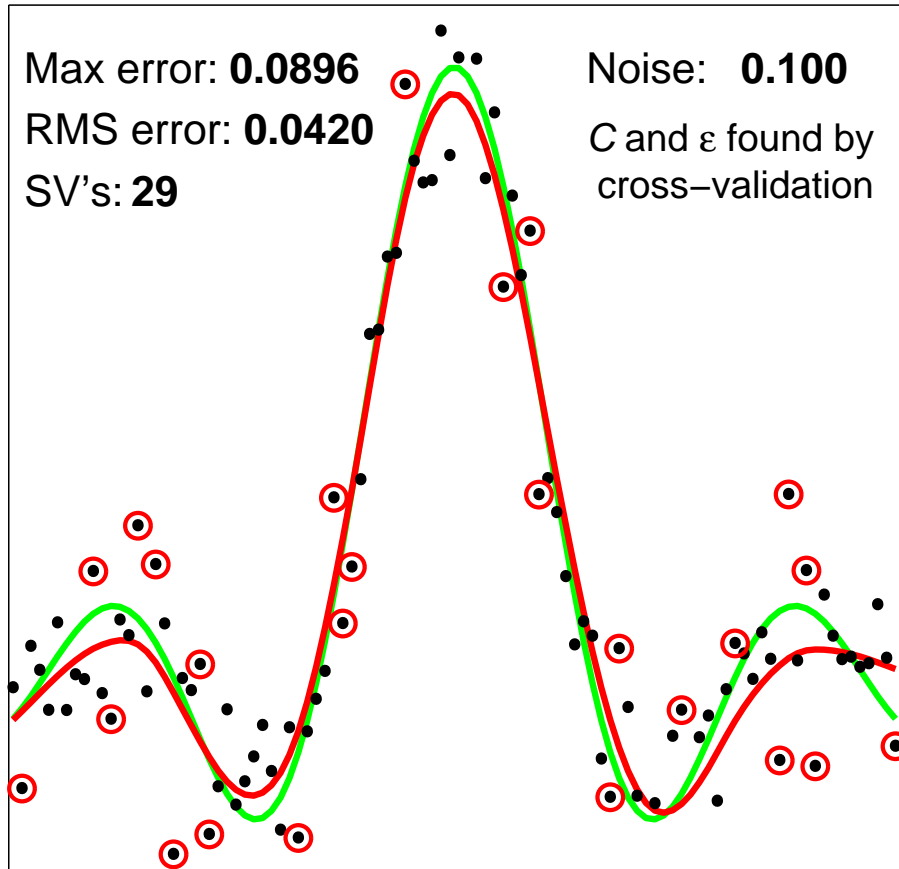
$\text{ML}(\mathbf{w})$

$\text{ML}(\alpha)$

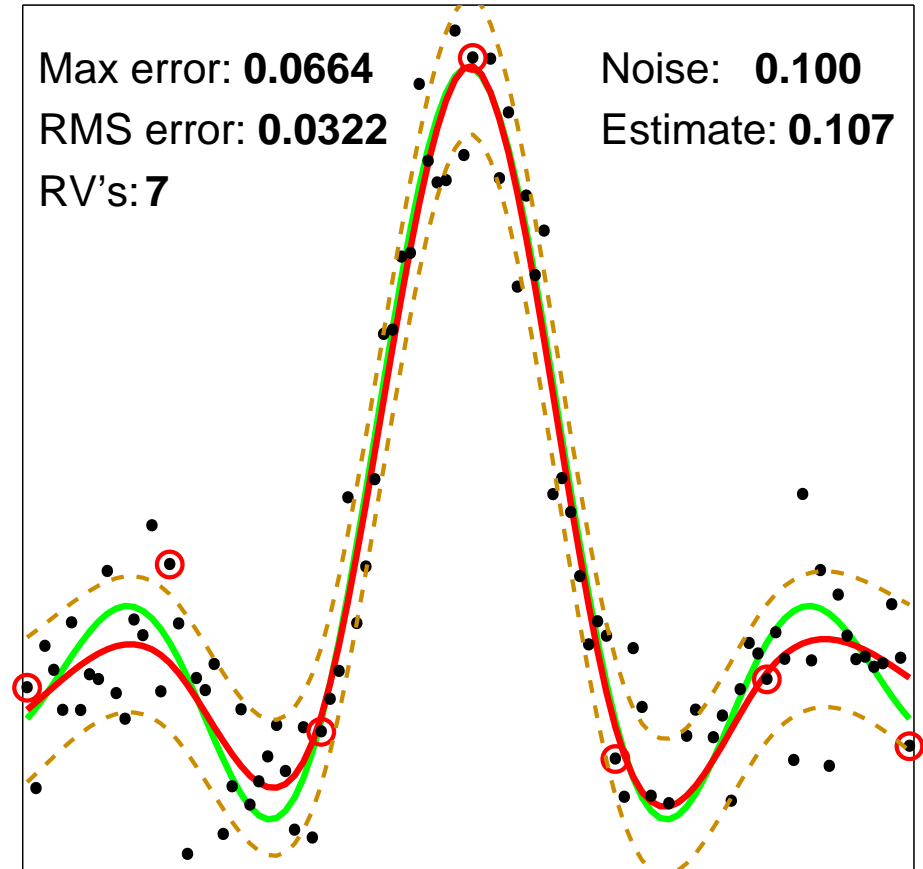
$$\mathbf{C} = \sum_{m=1}^M \alpha_m^{-1} \phi_m \phi_m^T + \sigma^2 \mathbf{I}$$



Support Vector Regression



Relevance Vector Regression



Sparse Bayesian Classification

- The likelihood is now *Bernoulli*, rather than Gaussian:

$$P(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N g\{y(\mathbf{x}_n; \mathbf{w})\}^{t_n} [1 - g\{y(\mathbf{x}_n; \mathbf{w})\}]^{1-t_n}$$

with $g(y) = 1/(1 + e^{-y})$ the ‘logistic sigmoid’ function

- Note: no noise variance σ^2
- Same sparse prior $p(\mathbf{w}|\alpha)$ as for regression
- Unlike for regression and the Gaussian likelihood model, $p(\mathbf{w}|\mathbf{t}, \alpha)$ cannot be obtained analytically, we so utilise a Laplace (Gaussian) approximation

Classification: Gaussian Posterior Approximation

- For the current values of α , find the posterior mode \mathbf{w}_{MP} via optimisation

- Compute the Hessian at \mathbf{w}_{MP} :

$$\nabla_{\mathbf{w}} \nabla_{\mathbf{w}} \log p(\mathbf{w}|\mathbf{t}, \alpha)|_{\mathbf{w}_{\text{MP}}} = -(\Phi^T \mathbf{B} \Phi + \mathbf{A})$$

where $\mathbf{B}_{nn} = g\{y(\mathbf{x}_n; \mathbf{w}_{\text{MP}})\} [1 - g\{y(\mathbf{x}_n; \mathbf{w}_{\text{MP}})\}]$

- Negate and invert to give the covariance Σ for a Gaussian approximation
 $p(\mathbf{w}|\mathbf{t}, \alpha) \approx N(\mathbf{w}_{\text{MP}}, \Sigma)$

- Hyperparameters α are updated as before using the approximated μ and Σ

- Note that $p(\mathbf{w}|\mathbf{t}, \alpha)$ is log-concave

Support Vector Machines

- A state-of-the-art method for classification and regression
- Given data set comprising N input vectors \mathbf{x}_n , model has form

$$y(\mathbf{x}; \mathbf{w}) = \sum_{n=1}^N w_n K(\mathbf{x}, \mathbf{x}_n) + w_0$$

- As many *kernel* functions $K(\cdot, \mathbf{x}_n)$ as examples $\Rightarrow M = N + 1$ parameters
- Support vector learning: minimise objective function of form:

$$\hat{E}(\mathbf{w}) = E_{\mathcal{D}}(\mathbf{w}) - \lambda \times (\text{size of margin})$$

- gives excellent accuracy (particularly in classification)
- as a side-effect, many w_n get set to zero — the model is *sparse*

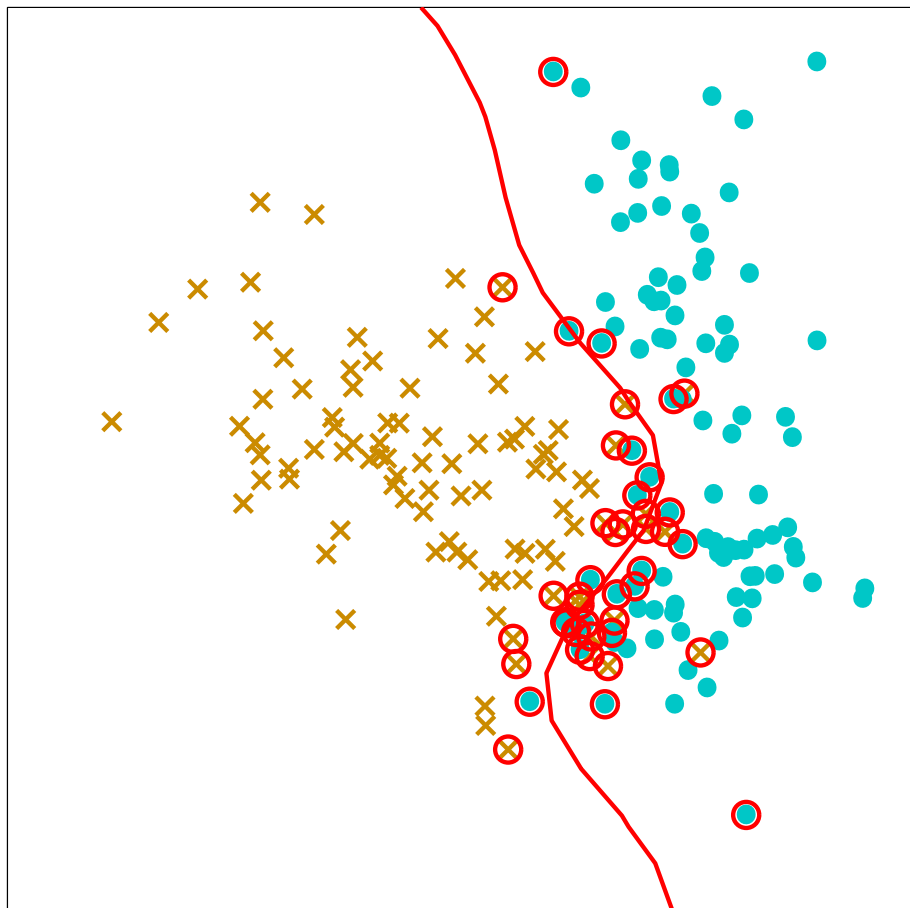
The “Relevance Vector Machine” (RVM)

- The RVM is simply a sparse Bayesian model utilising the same data-dependent kernel basis as the SVM:

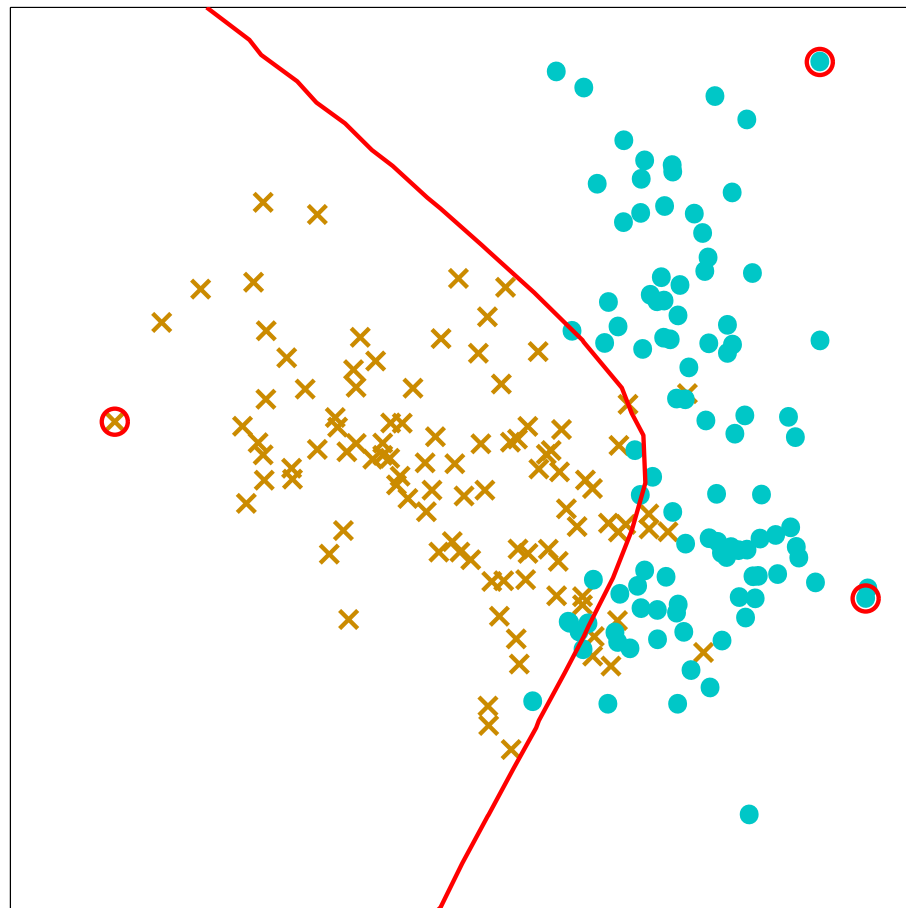
$$y(\mathbf{x}; \mathbf{w}) = \sum_{n=1}^N w_n K(\mathbf{x}, \mathbf{x}_n) + w_0$$

- Demonstration of RVM classification ...

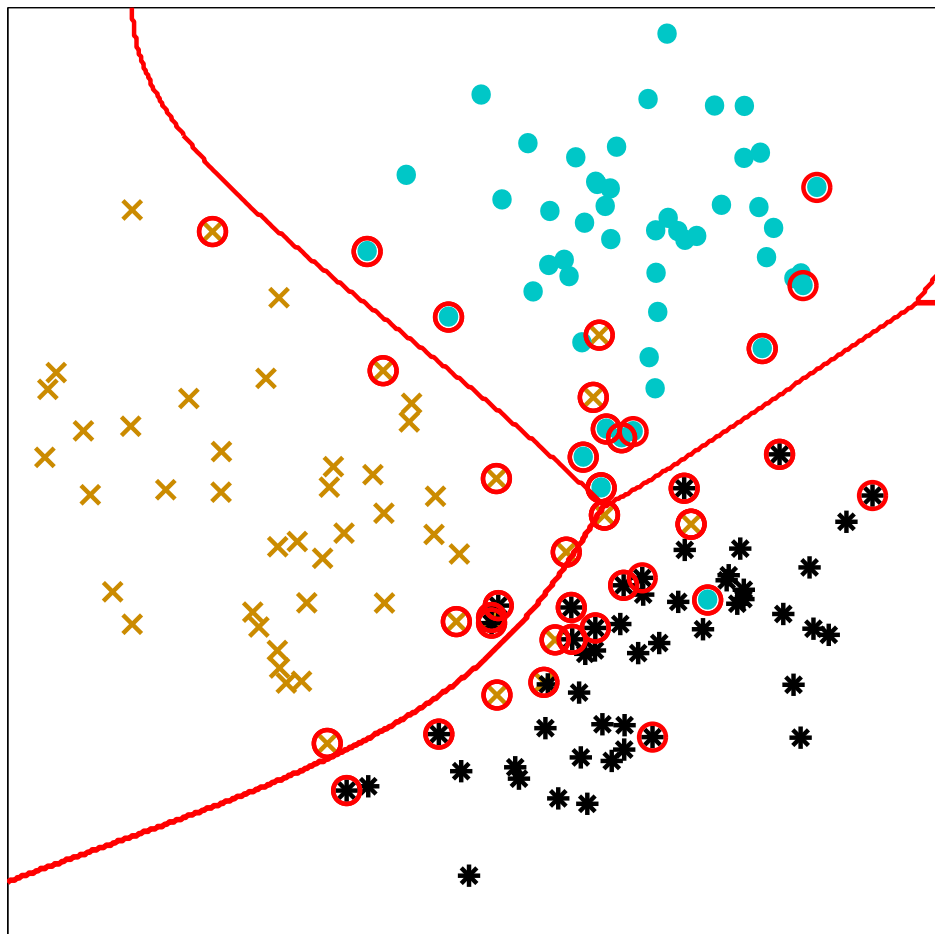
SVM: error=9.48% vectors=44



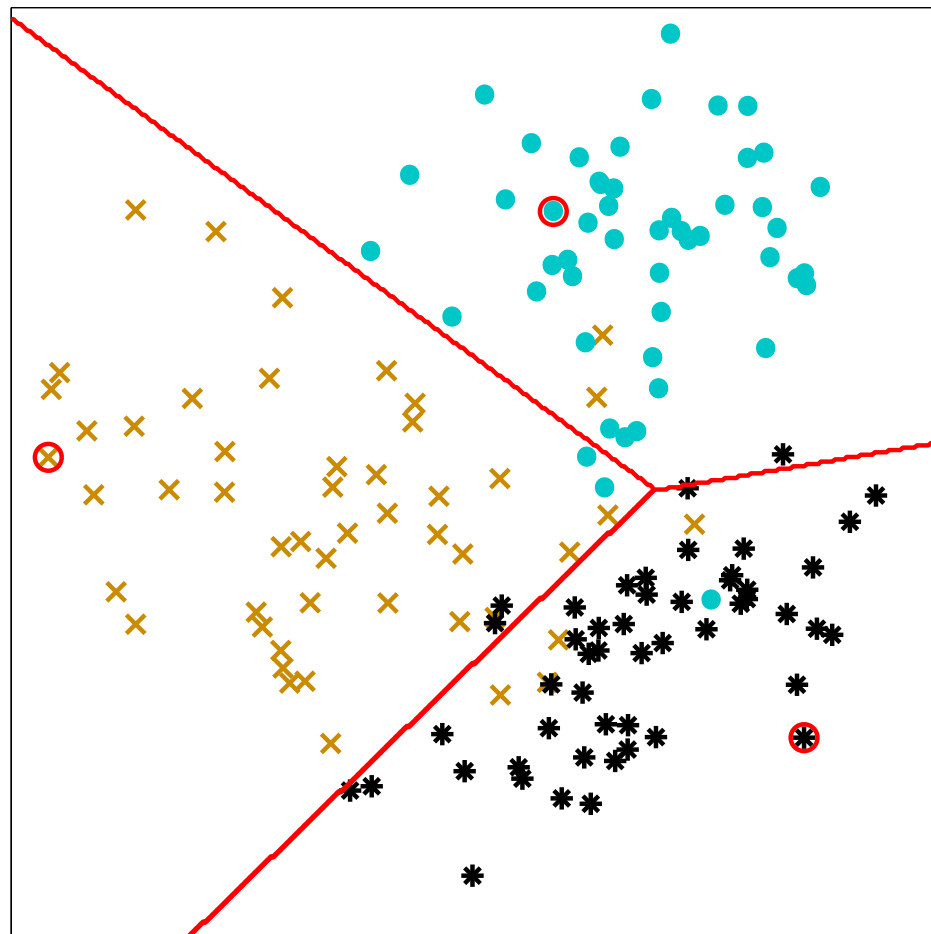
RVM: error=9.32% vectors=3



SVM: error=10.97% vectors=38



RVM: error=10.43% vectors=3



Comparison with the SVM

- General observations:

- RVM gives better generalisation in regression (?)
- SVM gives better generalisation in classification (?)
- RVM is much sparser (but the SVM is not designed to be sparse)

- There are other advantages of a Bayesian approach:

- no 'nuisance' parameters to set
- posterior probabilities in classification
- error bars in regression
- principled method for more than two classes
- not limited to Mercer kernels
- potential to estimate input scale parameters and compare kernels

Choosing the Kernel

- As in the SVM, we must choose the kernel and set any associated parameter(s)
- A Bayesian could compare alternative kernels by computing the fully marginalised probabilities of the data under candidate models. *e.g.*:

$$p(\mathbf{t}|\mathcal{K}_1) = \int p(\mathbf{t}|\alpha, \mathcal{K}_1) p(\alpha|\mathcal{K}_1) d\alpha$$

- We already know this integral isn't analytically tractable
- Approximation via sampling (as in Lecture 2) is not feasible for multiple α
- Deterministic approximations to this integral have proved inaccurate
- But $p(\mathbf{t}|\alpha_{\text{MP}}, \mathcal{K})$ is a 'reasonable' criterion for choosing kernels

Kernel 'Input Scale' Parameters

- For example, consider choice of η in $K(\mathbf{x}, \mathbf{x}_n) = \exp \left\{ -\eta \|\mathbf{x} - \mathbf{x}_n\|^2 \right\}$
- SVM: cross-validation can be used to set scale parameter η
- RVM: we can *optimise* the marginal likelihood function with respect to η
- Furthermore, we can optimise multiple scale parameters η , one for each of the d input dimensions:

$$K(\mathbf{x}, \mathbf{x}_n) = \exp \left\{ - \sum_{k=1}^d \eta_k (x_k - x_{nk})^2 \right\}$$

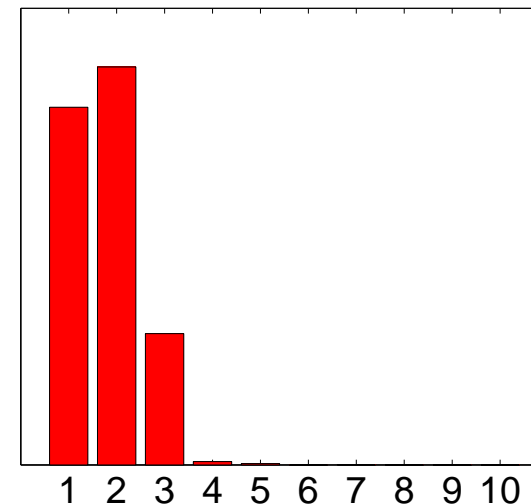
- Implementing sparsity of *input* variables (*q.v.* Gaussian process models)

Impact on Regression Benchmarks

- η -RVR: optimisation over both α and η

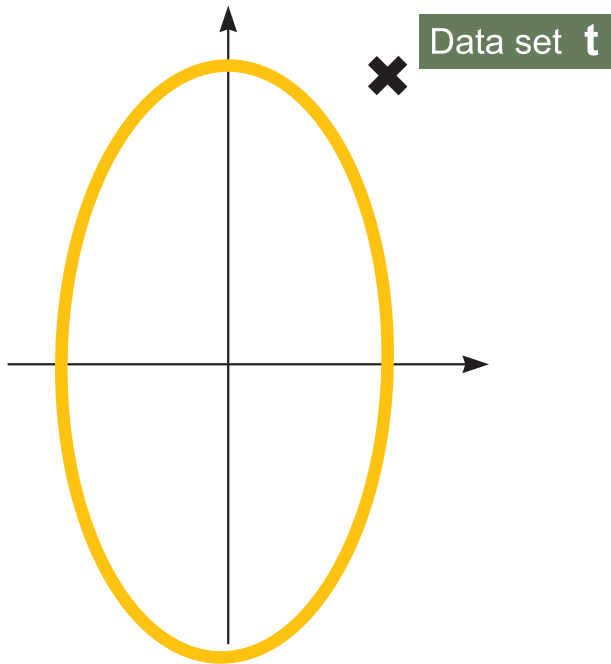
Dataset	Test error			# kernels		
	SVR	RVR	η -RVR	SVR	RVR	η -RVR
Friedman #1	2.92	2.80	0.27	116.6	59.4	11.5
Friedman #2	4140	3505	2593	110.3	6.9	3.9
Friedman #3	0.0202	0.0164	0.0119	106.5	11.5	6.4

- Friedman #1: 10-dimensional input space, but function depends only on variables 1–5. Final η -values shown on right:

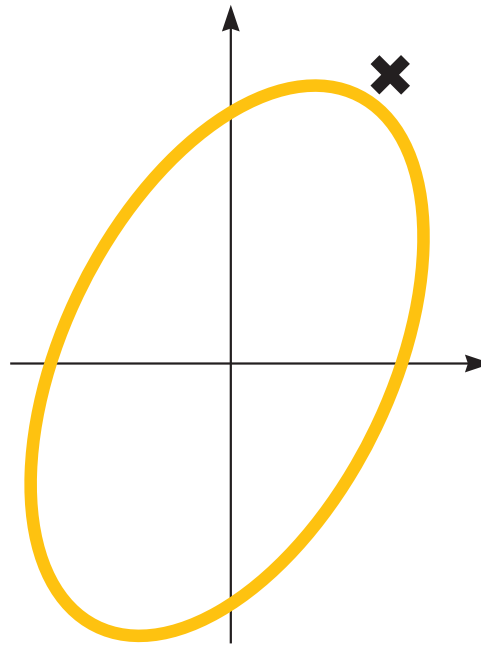


Ockham's Razor and Kernel Width

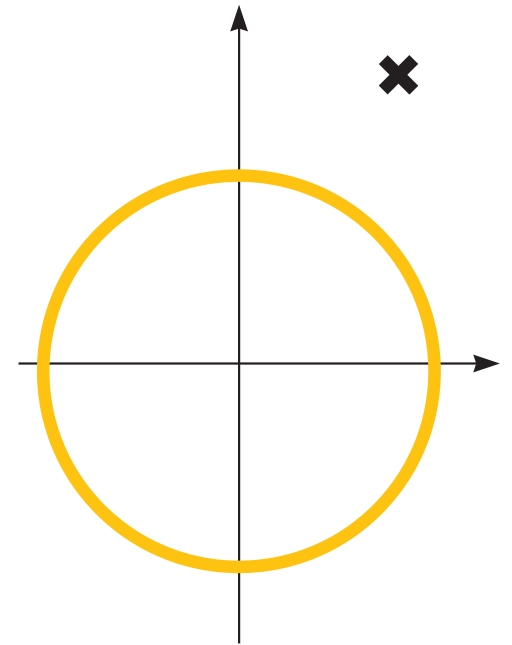
Observation Space \mathbf{R}^N



small basis function width



intermediate



large

Summary

- Experience appears to show that the sparse Bayesian learning procedure works highly effectively:
 - generalisation performance is typically very good
 - models are typically extremely (near-optimally?) sparse
- Challenge: obtain a reliable approximation to integrating out α (for model comparison)
- In Lecture 4, we'll look at:
 - Properties of the marginal likelihood function
 - An efficient algorithm for optimising α
 - Extensions and applications of sparse Bayesian models

Regression Performance Illustration

Data set	N	d	<u>errors</u>		<u>vectors</u>	
			SVM	RVM	SVM	RVM
Sinc (Gaussian noise)	100	1	0.378	0.326	45.2	6.7
Sinc (Uniform noise)	100	1	0.215	0.187	44.3	7.0
Friedman #1	240	10	2.92	2.80	116.6	59.4
Friedman #2	240	4	4140	3505	110.3	6.9
Friedman #3	240	4	0.0202	0.0164	106.5	11.5
Boston Housing	481	13	8.04	7.46	142.8	39.0
Normalised Mean			1.00	0.86	1.00	0.15

Classification Performance Illustration

Data set	N	d	<u>errors</u>		<u>vectors</u>	
			SVM	RVM	SVM	RVM
Pima Diabetes	200	8	20.1%	19.6%	109	4
U.S.P.S.	7291	256	4.4%	5.1%	2540	316
Banana	400	2	10.9%	10.8%	135.2	11.4
Breast Cancer	200	9	26.9%	29.9%	116.7	6.3
Titanic	150	3	22.1%	23.0%	93.7	65.3
Waveform	400	21	10.3%	10.9%	146.4	14.6
German	700	20	22.6%	22.2%	411.2	12.5
Image	1300	18	3.0%	3.9 %	166.6	34.6
Normalised Mean			1.00	1.08	1.00	0.17