

Bayesian Inference

Mohammad Emtiyaz Khan
AIP (RIKEN), Tokyo

<http://emtiyaz.github.io>

emtiyaz.khan@riken.jp

June 28-29, 2018



©Mohammad Emtiyaz Khan 2018

Contents

1	Introduction	2
2	Regression and Classification	3
3	Maximum Likelihood	4
4	Maximum A Posteriori (MAP)	6
5	The Posterior Distribution	6
6	Bayesian Linear Regression	11
7	Bayesian Logistic Regression	13
8	Deep Neural Networks	16
9	Gaussian Processes	17
10	Benefits of Bayesian Inference	19
11	Summary	24
	List of concepts	26

1 Introduction

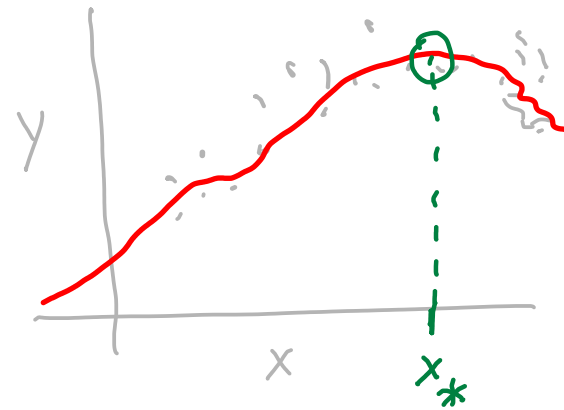
We will consider Bayesian inference in the context of supervised learning. Bayesian inference involves computation of posterior distribution, which is fundamentally different from the maximum-likelihood principle. We will demonstrate this on four models: linear regression, logistic regression, Neural networks, and Gaussian process.

By using the posterior distribution, Bayesian inference can reduce overfitting, represent uncertainty, and perform model selection.

Unfortunately, Bayesian inference involves a difficult integral which involves computing an average over all possible explanations of the data. We will learn about some of the reasons behind this difficulty. Approximate Bayesian inference addresses this problem by finding approximations to the integral

2 Regression and Classification

Regression/classification is to relate input variables to the output variable, to predict outputs for new inputs and/or to understand the effect of the input on the output.



Dataset for regression

The data, denoted by \mathcal{D} , consist of pairs (\mathbf{x}_n, y_n) , where \mathbf{x}_n is a vector of D inputs and y_n is the n 'th output. Number of pairs N is the data-size and D is the dimensionality.

$$y_n, \underline{x}_n$$
$$y_n \approx f(\underline{x}_n)$$

Prediction

In prediction, we wish to predict the output for a new input vector, i.e., find a regression function that approximates the output “well enough” given inputs.

$$y_n \approx f_{\mathbf{w}}(\mathbf{x}_n), \text{ for all } n$$

where \mathbf{w} is the parameter of the regression model.

3 Maximum Likelihood

Assume y_n to be independent samples from an *exponential-family distribution*, whose *expectation parameter* is equal to $f_w(\mathbf{x}_n)$:

$$p(\mathcal{D}|\mathbf{w}) := \prod_{n=1}^N p(y_n | f_w(\mathbf{x}_n)).$$

The function $p(\mathcal{D}|\mathbf{w})$ is the **likelihood**, which can be maximized to obtain a “good enough” \mathbf{w} ,

$$\mathcal{L}_{ML}(\mathbf{w}) := \log p(\mathcal{D}|\mathbf{w}).$$

This is known as the **maximum-likelihood estimation**.

$$P(y_n | f_w(\mathbf{x}_n))$$

$$\prod_{n=1}^N \mathcal{N}(y_n | f_w(\mathbf{x}_n), 1)$$

log ↓

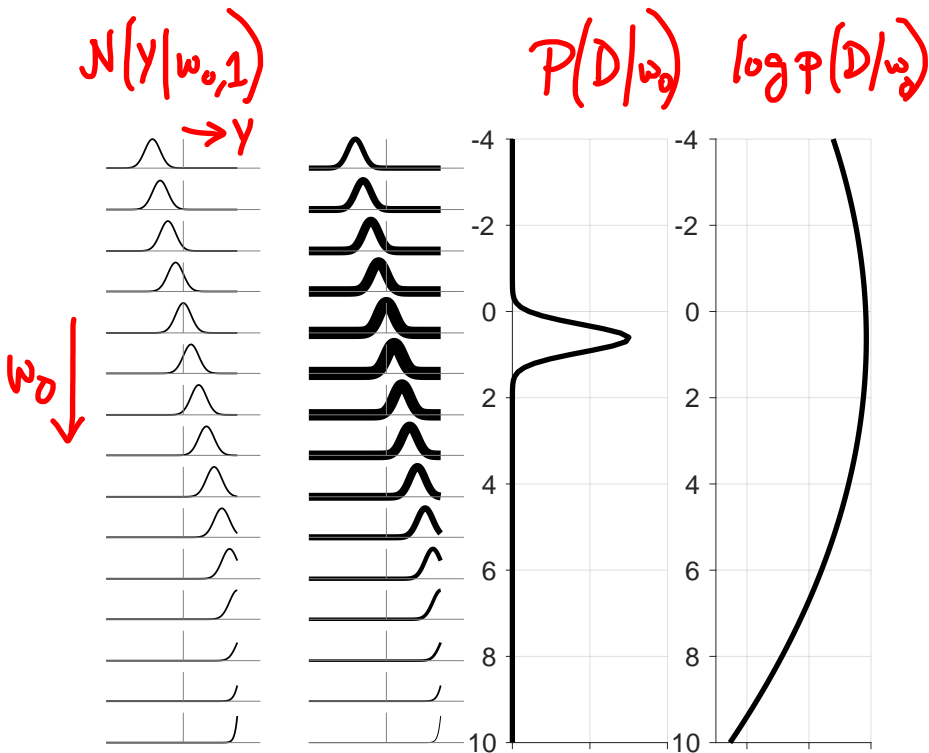
$$\sum_{n=1}^N (y_n - f_w(\mathbf{x}_n))^2$$

1 parameter model

$$y_n \approx w_0$$

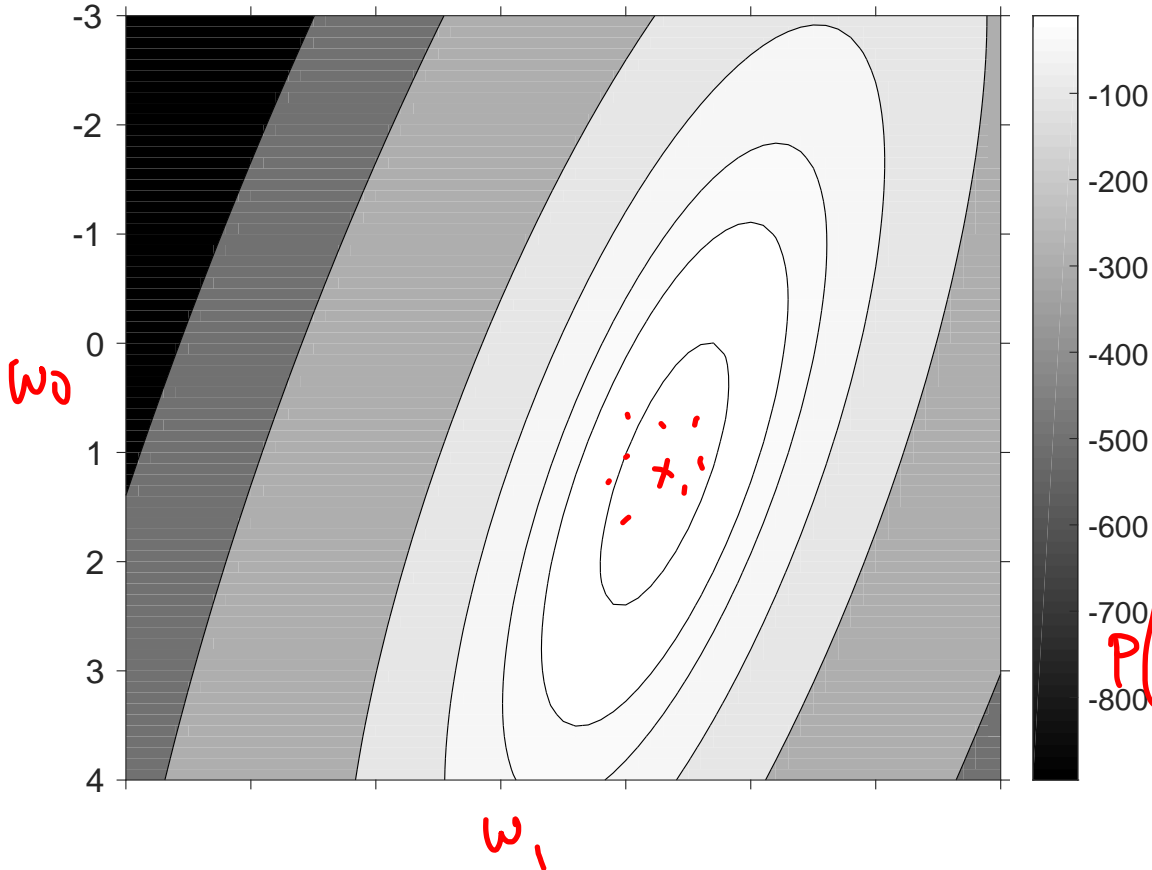
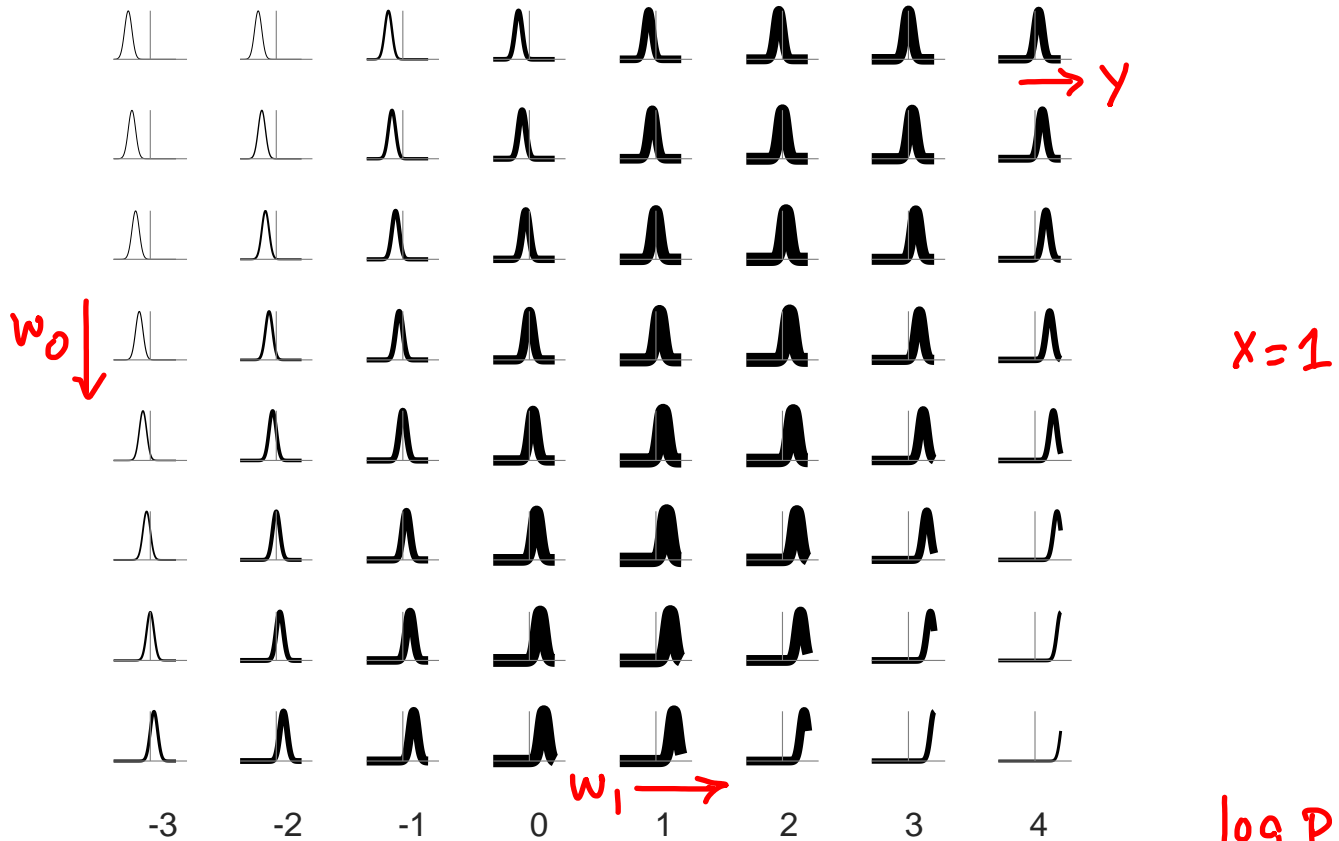
$$\mathcal{N}(y_n | w_0, 1)$$

$$\sum_{n=1}^N \log \mathcal{N}(y_n | w_0, 1)$$

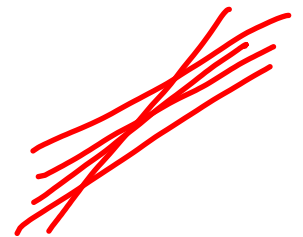


Two-parameter model.

$$y = w_0 + w_1 x$$



$$\log P(D|w_0, w_1)$$



$$P(w|D) = \frac{P(D|w)}{\int P(D|w) dw}$$

4 Maximum A Posteriori (MAP)

To avoid overfitting, we can use a regularizer $\log p(\mathbf{w})$ to perform maximum a posteriori estimation,

$$R(\omega) = \lambda \omega^T \omega$$

$$= \log \mathcal{N}(\omega | 0, I/\lambda)$$

$$\mathcal{L}_{MAP}(\mathbf{w}) := \log p(\mathcal{D}|\mathbf{w}) + \log p(\mathbf{w}).$$

Note that not all regularizers correspond to a probability distribution.

We can view $p(\mathbf{w})$ as a prior distribution to get the joint distribution:

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}),$$

$$= p(\mathbf{w}|\mathcal{D})p(\mathcal{D}).$$

$$P(A, B) = P(A|B)P(B)$$

$$= P(B|A)P(A)$$

$$\frac{P(A|B)P(B)}{P(A)} = \frac{P(A, B)}{P(A)} = P(B|A)$$

5 The Posterior Distribution

The posterior distribution is defined using the Bayes' rule:

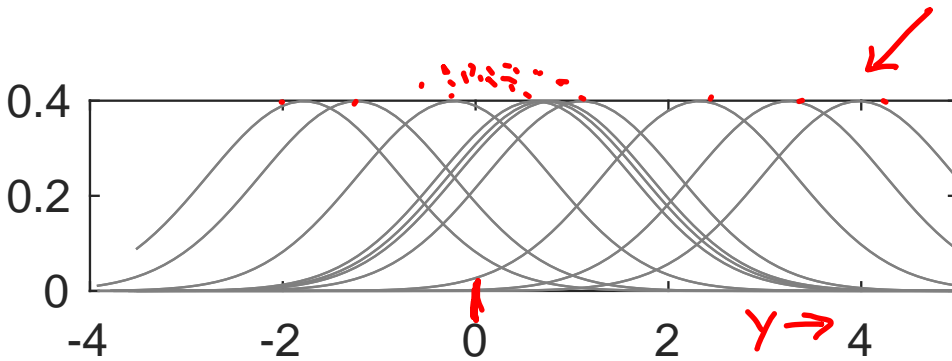
$$\underline{p(\mathbf{w}|\mathcal{D})} = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{\int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}}.$$

$$\frac{P(A|B)P(B)}{\int P(A|B)P(B)dB} = P(A)$$

The integral $p(\mathcal{D})$ is the normalizing constant, also known as, the marginal likelihood. Without it, we do not know the true spread of the distribution, which gives us a notion of uncertainty or confidence.

Example: One-Parameter Model

$$y_n \approx \omega_0$$

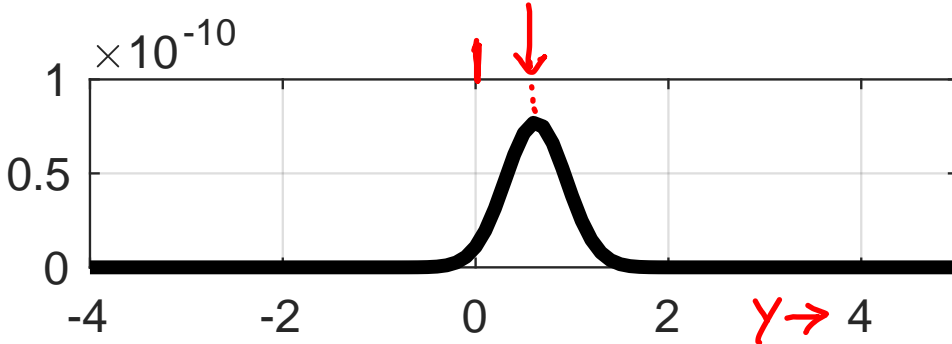


$$\mathcal{N}(y_n | \omega_0, 1)$$

lik

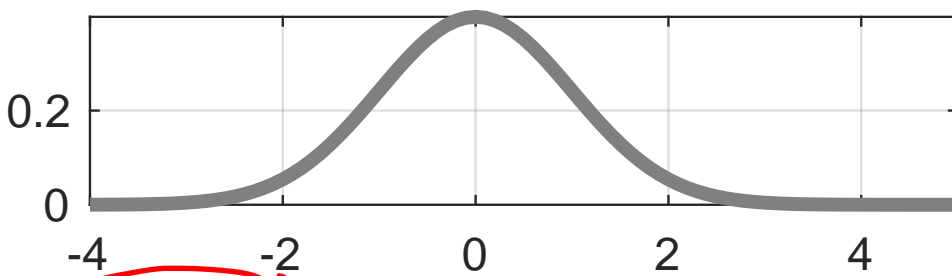
$$\prod_{n=1}^N \mathcal{N}(y_n | \omega_0, 1)$$

Prior



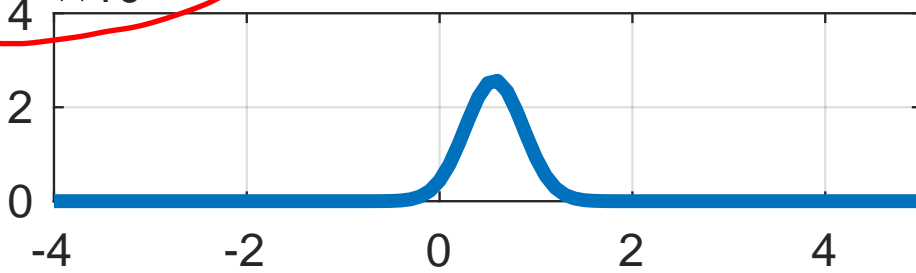
$$\mathcal{N}(\omega_0 | 0, 1)$$

$$\text{lik}(\omega_0) \text{Prior}(\omega_0)$$

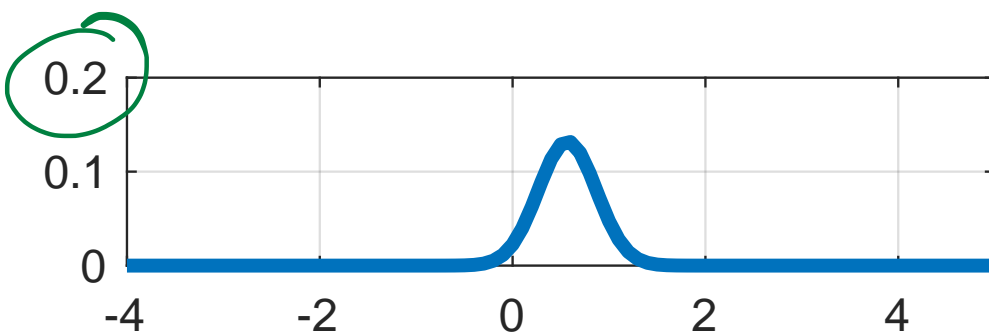


$$\sum_{\omega_0} \text{lik}(\omega_0) \text{Prior}(\omega_0)$$

$$4 \times 10^{-11}$$



$$P(\omega_0 | D)$$



Posterior of a One-Parameter Model

$$\mathcal{N}(z|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\frac{(z-\mu)^2}{\sigma^2}}$$

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{\int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

lik Prior

$$p(w_0|\mathcal{D}) = \frac{\left[\prod_{n=1}^N \mathcal{N}(y_n|w_0, 1) \right] \mathcal{N}(w_0|0, 1)}{\int \left[\prod_{n=1}^N \mathcal{N}(y_n|w_0, 1) \right] \mathcal{N}(w_0|0, 1)dw_0}$$

lik Prior

Question: Derive the posterior distribution and the marginal likelihood.

$$w_{MAP} = \frac{1}{N+1} \sum y_n$$

$$\mu \triangleq \frac{1}{N+1} \sum y_n$$

$$\sigma^2 \triangleq \frac{1}{N+1}$$

$$\left[\prod_{n=1}^N e^{-\frac{1}{2}(y_n - w_0)^2} \right] e^{-\frac{1}{2}w_0^2}$$

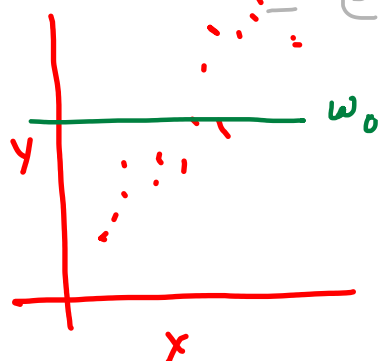
$$= e^{-\frac{1}{2} \left[\sum_{n=1}^N (y_n - w_0)^2 + w_0^2 \right]}$$

$$= e^{-\frac{1}{2} \left[\sum_n y_n^2 + Nw_0^2 - 2\left(\sum_n y_n\right)w_0 + w_0^2 \right]}$$

$$= e^{-\frac{1}{2} \left[(N+1)w_0^2 - 2\left(\sum_n y_n\right)w_0 + \sum_n y_n^2 \right]}$$

$$= e^{-\frac{1}{2} \left[(N+1) \left[w_0^2 - 2\left(\frac{1}{N+1} \sum_n y_n\right)w_0 + \mu^2 \right] + \sum_n y_n^2 - \mu^2 \right]}$$

$$= \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{1}{2} \frac{(N+1)(w_0 - \mu)^2}{\sigma^2}} \times \sqrt{2\pi}\sigma^2 e^{-\frac{1}{2} \left(\sum_n y_n^2 - \frac{1}{N+1} \left(\sum_n y_n \right)^2 \right)}$$



$$p(\mathcal{D}) =$$

$$p(w_0|\mathcal{D}) = \mathcal{N}(w_0|\mu, \sigma^2) \mathcal{N}\left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \middle| \cdot, \cdot\right)$$

P(w₀|D) P(D)

Example: Two-Parameter Model

$$y_n \approx w_0 + w_1 x_n$$

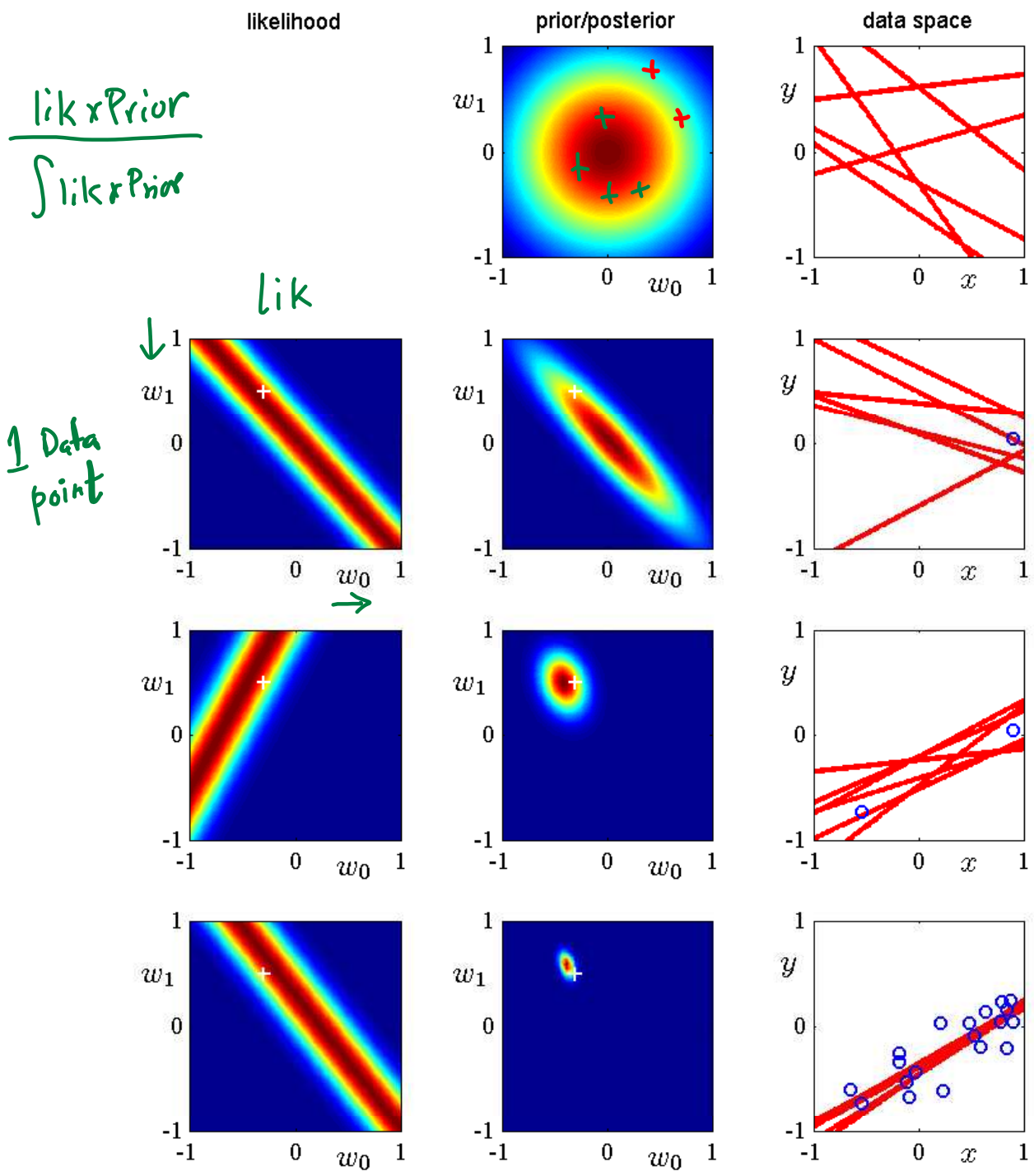


Figure taken from [Bishop, 2006]

Example: Two-Parameter Model

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{\int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}},$$

$$p(\mathbf{w}|\mathcal{D}) = \frac{\left[\prod_{n=1}^N \mathcal{N}(y_n|w_0 + w_1x_n, 1) \right] \mathcal{N}(\mathbf{w}|0, \mathbf{I})}{\int \left[\prod_{n=1}^N \mathcal{N}(y_n|w_0 + w_1x_n, 1) \right] \mathcal{N}(\mathbf{w}|0, \mathbf{I})d\mathbf{w}}.$$

Question: Derive the posterior distribution and the marginal likelihood.

$$[w_0 \ w_1] \begin{bmatrix} 1 \\ x_n \end{bmatrix}$$

$$-\frac{1}{2} \left[\sum_{n=1}^N (y_n - \omega^T \tilde{x}_n)^T (y_n - \omega^T \tilde{x}_n) + \omega^T \omega \right]$$

$$\sum_n \left[y_n y_n - 2 y_n x_n^T \omega + \omega^T x_n x_n^T \omega \right] + \omega^T \omega$$

$$\omega^T \underbrace{\left(\mathbf{I} + \mathbf{X}^T \mathbf{X} \right)}_{\mathbf{V}^{-1}} \omega - 2 \underbrace{\mathbf{y}^T \mathbf{X}}_{\mathbf{V}^{-1} \mathbf{m}} \omega + \text{const}$$

$$p(\mathcal{D}) =$$

$$p(w_0, w_1|\mathcal{D}) =$$

6 Bayesian Linear Regression

Consider $f_{\mathbf{w}}(\mathbf{x}_n) := \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n)$,
 (where $\boldsymbol{\phi}(\cdot)$ is a nonlinear function)
 with the following joint distribution,

← Basis-function expansion
 $[\omega_0 \ \omega_1 \ \dots \ \omega_D]^\top$
 $\begin{bmatrix} 1 \\ \phi_1(x_n) \\ \phi_2(x_n) \\ \vdots \\ \phi_D(x_n) \end{bmatrix}$

$$p(\mathcal{D}, \mathbf{w}) = \left[\prod_{n=1}^N \mathcal{N}(y_n | \overset{\text{nonlinear}}{f_{\mathbf{w}}(\mathbf{x}_n)}, 1) \right] \mathcal{N}(\mathbf{w} | 0, \lambda \mathbf{I}).$$

Question: Derive the posterior distribution and the marginal likelihood.

$$\mathbf{V}^{-1} = (\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})$$

$$\mathbf{V}^{-1} \mathbf{m} = \boldsymbol{\Phi}^\top \mathbf{y}$$

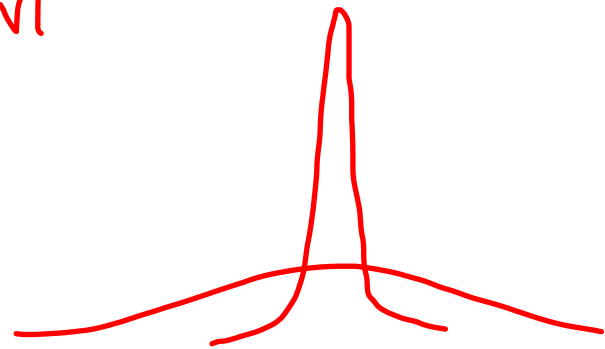
$$p(\mathbf{w} | \mathcal{D}) = \mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{V}),$$

$$\mathbf{m} = (\lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

$$\mathbf{V} = (\lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$$

$$p(\mathcal{D}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mathbf{I} / \lambda)$$

$$= \frac{1}{(2\pi)^{D/2} |\mathbf{V}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{w} - \mathbf{m})^\top \mathbf{V}^{-1} (\mathbf{w} - \mathbf{m})\right)$$



Question: When is the computation of the posterior distribution difficult?

The predictive distribution is,

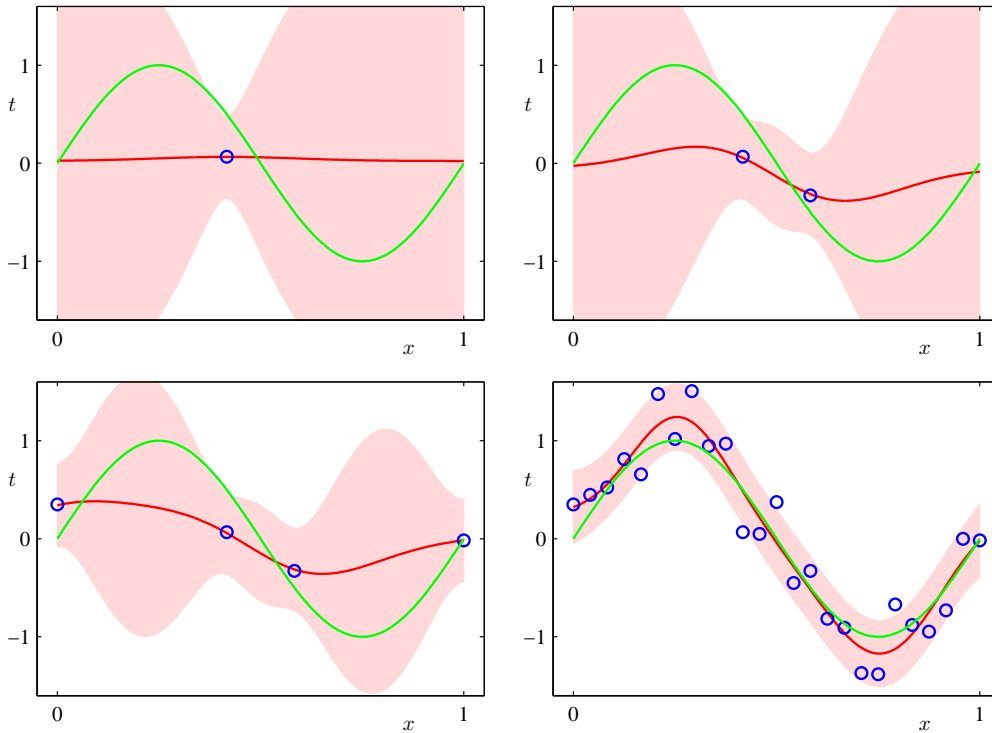
$$= \int \mathcal{P}(y_* | f_{\mathbf{w}}(x_*)) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

$$\underline{\underline{p(y_* | \mathbf{x}_*, \mathcal{D})}} = \int \mathcal{N}(y_* | f_{\mathbf{w}}(\mathbf{x}_*), 1) p(\mathbf{w} | \mathcal{D}) d\mathbf{w},$$

$$= \mathcal{N}(y_* | \mathbf{m}^\top \boldsymbol{\phi}(\mathbf{x}_*), 1 + \boldsymbol{\phi}(\mathbf{x}_*)^\top \mathbf{V} \boldsymbol{\phi}(\mathbf{x}_*)).$$

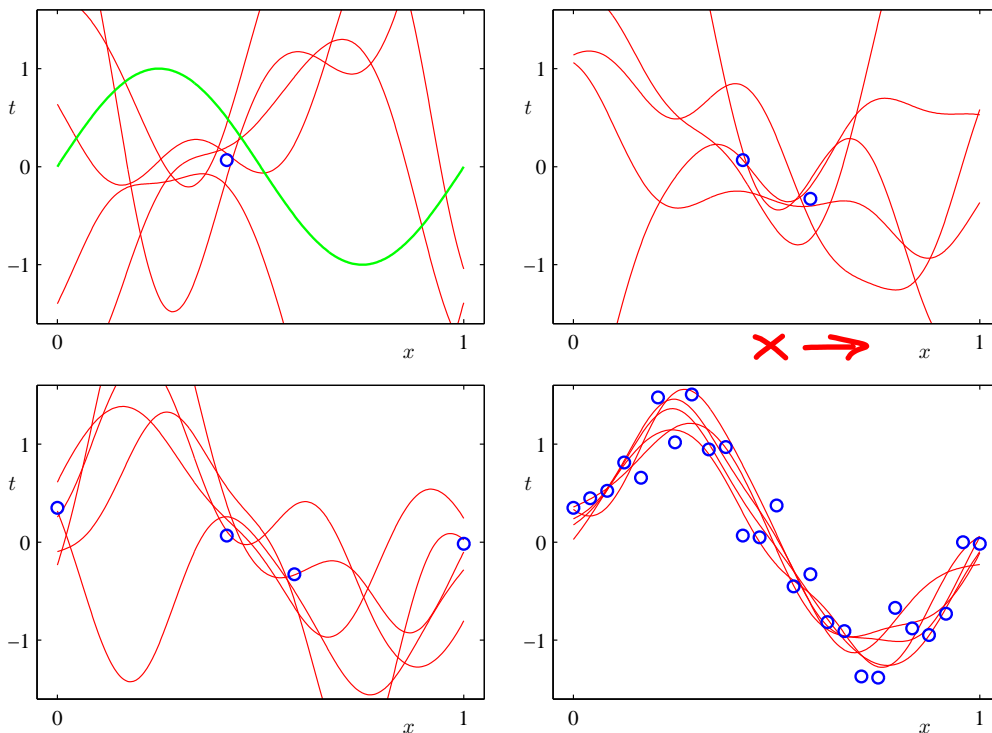
Example of Predictive Distribution

From [Bishop, 2006] Figure 3.8.



Sampled Prediction Functions

From [Bishop, 2006] Figure 3.9.



$$f_w(x)$$

$\times \rightarrow$

7 Bayesian Logistic Regression

When $y_n \in \{0, 1\}$, we can use a Bernoulli distribution,

$$p(\mathcal{D}, \mathbf{w}) = \left[\prod_{n=1}^N \text{Ber}(y_n | f_{\mathbf{w}}(\mathbf{x}_n)) \right] \mathcal{N}(\mathbf{w} | 0, \lambda \mathbf{I}),$$

$$p(y_n = 1 | f_{\mathbf{w}}(\mathbf{x}_n)) = \frac{1}{1 + e^{\mathbf{w}^\top \mathbf{x}_n}}. \quad \leftarrow \quad w_1 x_{n1} + w_2 x_{n2}$$

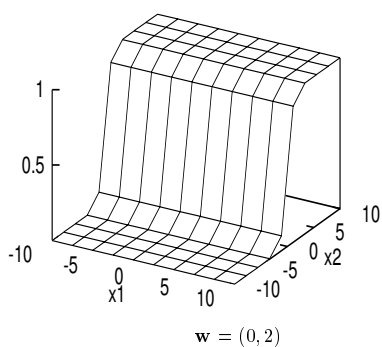
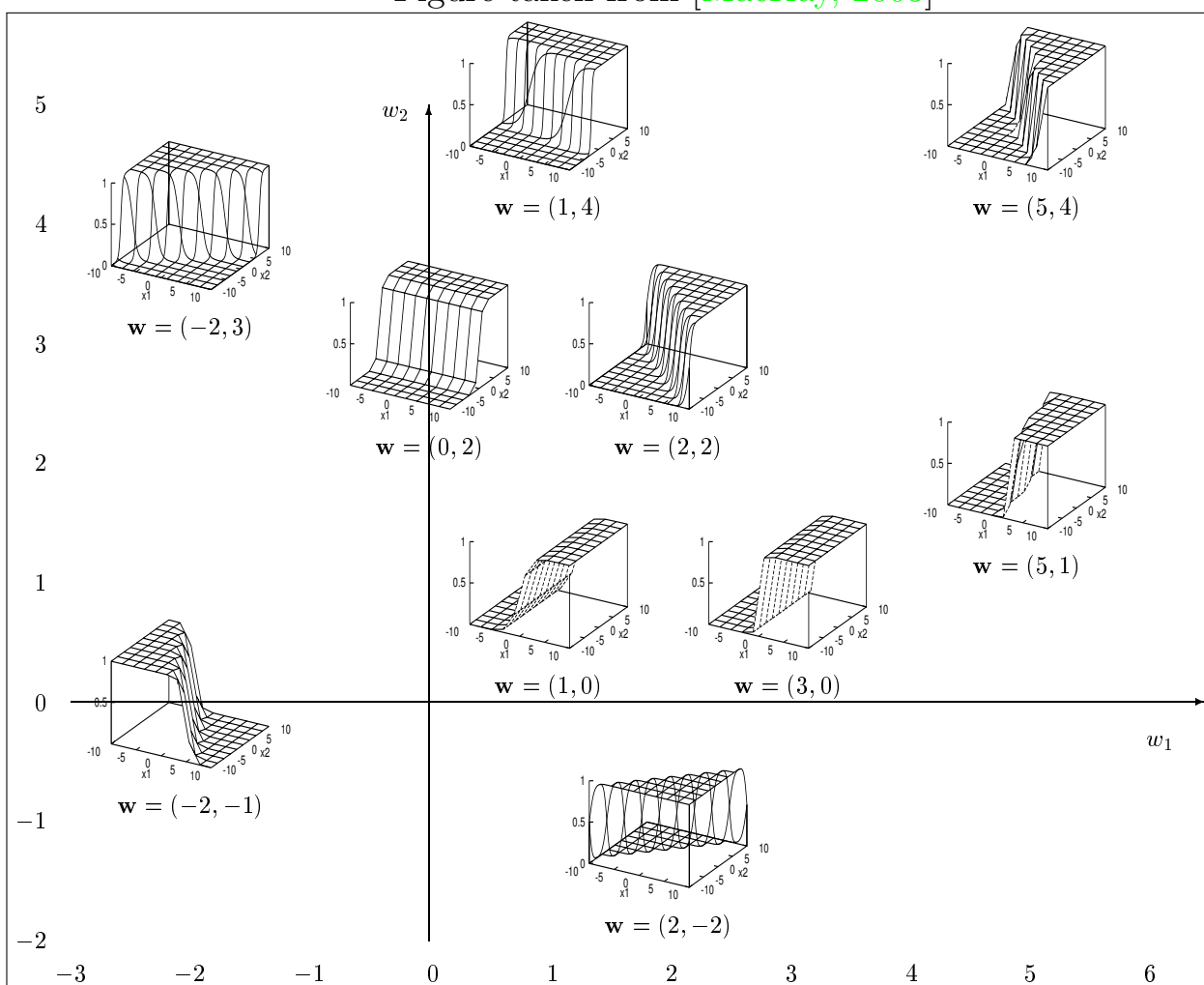
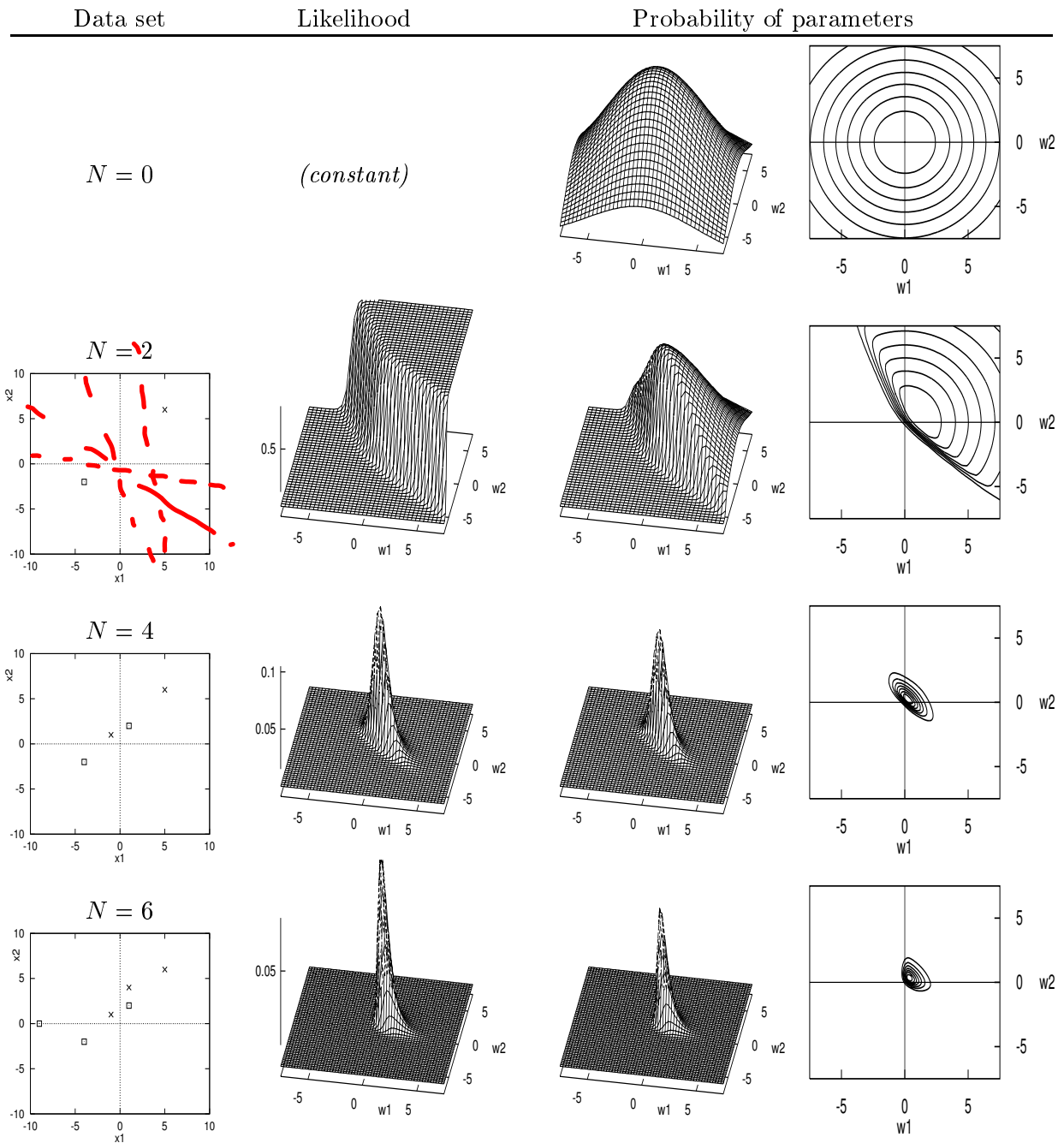
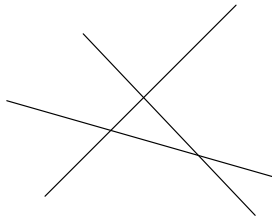


Figure taken from [MacKay, 2003]



2D Example



Posterior Distribution

The posterior distribution is,

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

The marginal-likelihood is,

$$p(\mathcal{D}) = \int \left\{ \prod_{n=1}^N \left[\frac{1}{1+e^{w^T x_n}} \right]^{y_n} \left[1 - \frac{1}{1+e^{w^T x_n}} \right]^{1-y_n} \right\} \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}) d\mathbf{w}$$

There is no closed-form solution. This is because the likelihood is not **conjugate** to the Gaussian prior. That is, with respect to \mathbf{w} , the likelihood cannot be expressed in the same form as the prior.

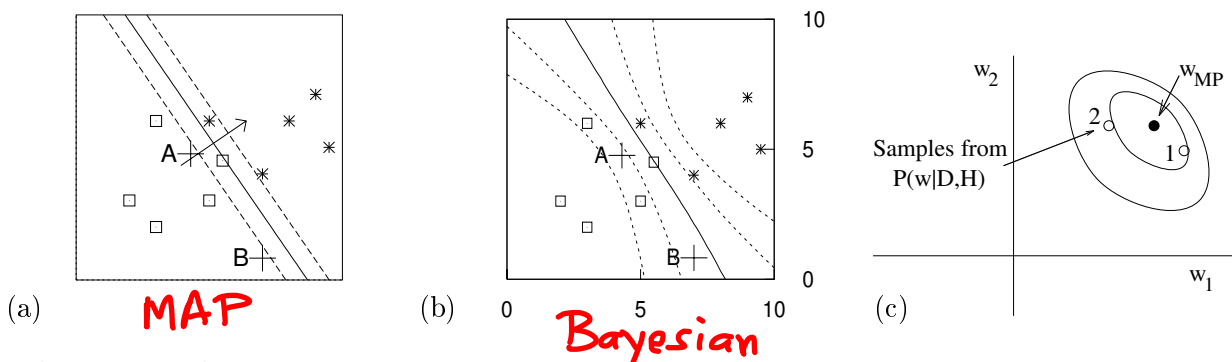
$$-\frac{1}{2} \mathbf{w}^T \mathbf{w}$$

Predictive Distribution

The predictive distribution is intractable too,

$$p(y_* | \mathbf{x}_*, \mathcal{D}) = \int \text{Ber}(y_* | f_w(\mathbf{x}_*)) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}.$$

Figure taken from [MacKay, 2003]

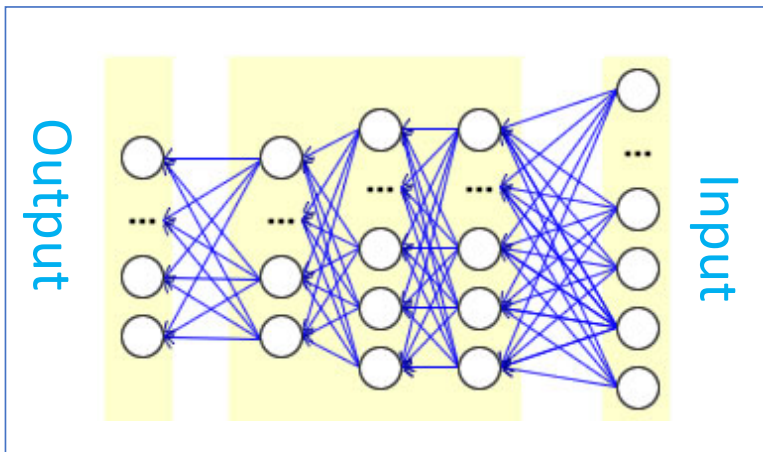


The prediction uncertainty is useful to avoid “overconfident” decision boundary found by a MAP method.

8 Deep Neural Networks

These difficult become much worse with deep models.

$$f_w(\mathbf{x}_n) = f_1(\mathbf{w}_1 f_2(\mathbf{w}_2 \dots f_L(\mathbf{w}_L \mathbf{x}_n)))$$



Nonconvexity translates to a multi-modal posterior distribution, whose marginal-likelihood is a massive, intractable integral.

Three computational challenges:

- 1) Too many factors (large N)
- 2) Too many dimensions (large D)
- 3) Nonconjugacy

9 Gaussian Processes

Instead of assigning a Gaussian prior to \mathbf{w} , we can directly apply an [informative prior](#) on $f_w(\mathbf{x})$ using [Gaussian process](#), which is defined using a kernel matrix \mathbf{K}_θ with entries $k_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$ where θ are kernel hyperparameters, e.g.,

$$k_{ij} := \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j).$$

We can then directly sample N function values $\mathbf{f} = [f_1, f_2, \dots, f_N]$ from a [Gaussian-process prior](#):

$$p(\mathbf{f}|\mathbf{X}, \theta) := \mathcal{N}(\mathbf{f}|0, \mathbf{K}_\theta).$$

No need for \mathbf{w} ! Combining it with a likelihood, we get GP models.

A good text is [[Rasmussen and Williams, 2006](#)]

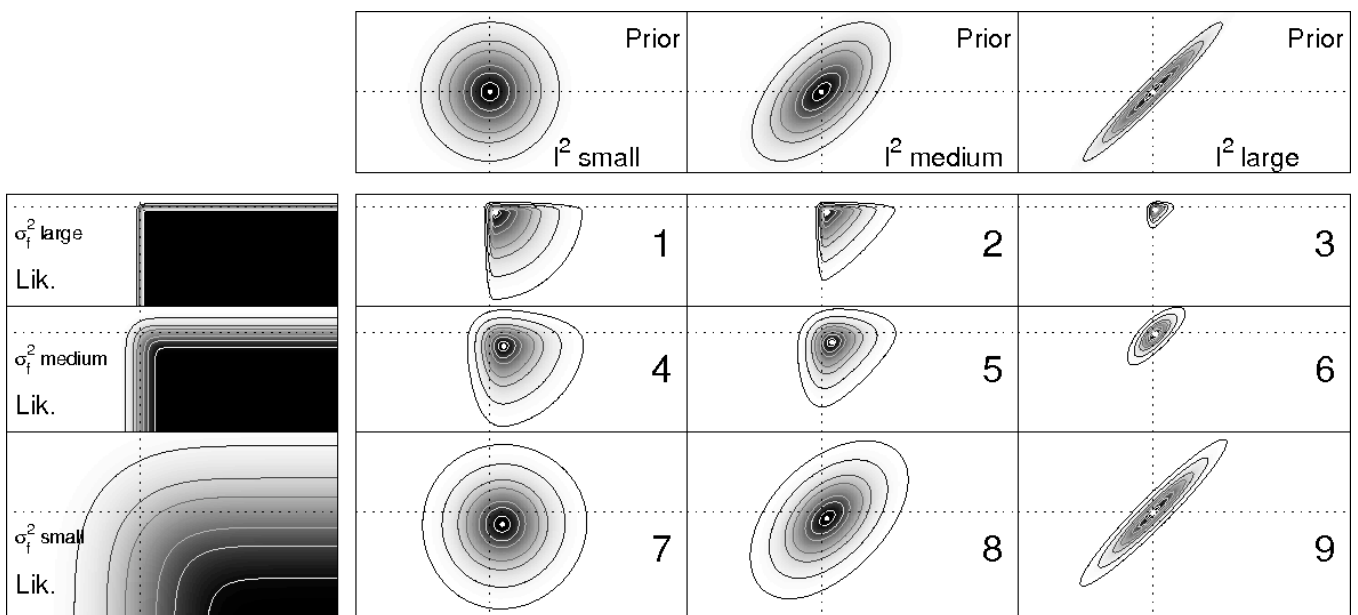
2D Example

Consider [square-exponential kernel](#):

$$k_{ij} := \sigma_f^2 \exp\left(-\frac{1}{2l}\|\mathbf{x}_i - \mathbf{x}_j\|^2\right),$$

along with a Bernoulli likelihood:

$$p(\mathcal{D}, \mathbf{w}) = \left[\prod_{n=1}^N \text{Ber}\left(y_n \mid \frac{1}{1 + \exp(f_n)}\right) \right] \mathcal{N}(\mathbf{f} | 0, \mathbf{K}_\theta).$$



Taken from [\[Nickisch and Rasmussen, 2008\]](#)

This is again a nonconjugate model, which results in an intractable integral. Could you think of a lower bound on the computation?

Hint: Use a Gaussian likelihood.

10 Benefits of Bayesian Inference

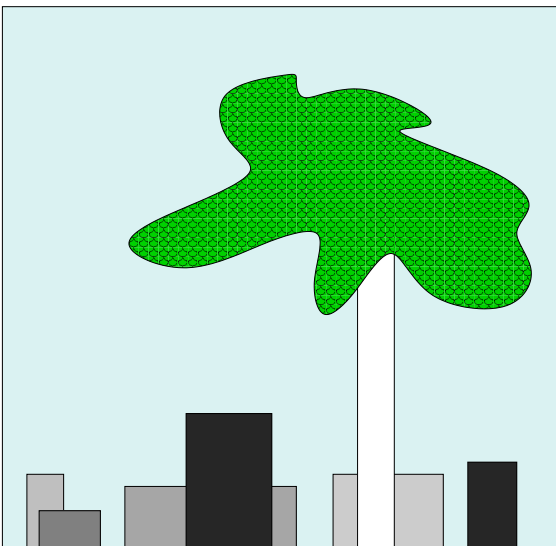
There are multiple benefits:

- 1) Posterior distribution gives an estimate of uncertainty.
- 2) Averaging with respect to it can reduce overfitting.
- 3) Posterior distribution enables data-generation.
- 4) Marginal likelihood enables model-selection.

We give more details about the last point now.

Occam's Razor

What is behind the tree?



From [MacKay, 2003].

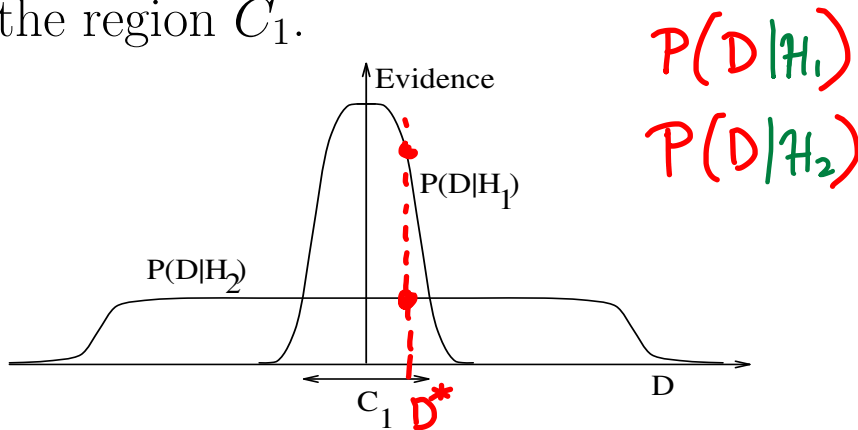
Simpler explanations are better – in absence of certainty. Bayesian inference naturally incorporates this principle.

Bayesian Razor

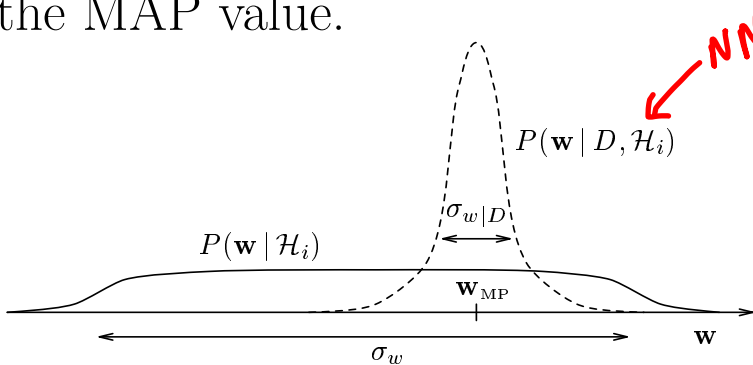
A simple model \mathcal{H}_1 can make only limited range of predictions well, while a complex models \mathcal{H}_2 can cover more range. However, this means that \mathcal{H}_2 does not predict the data set in region C_1 as strongly as \mathcal{H}_1 . Suppose the two models have equal prior probabilities, then \mathcal{H}_1 will be more probable than \mathcal{H}_2 in the region C_1 .

Suppose:

Linear model \mathcal{H}_1
 NN \mathcal{H}_2



In a similar fashion, at the parameter level the posterior distribution $p(\mathbf{w}|\mathcal{D}, \mathcal{H}_i)$ incorporates the Occam factor due to the “spread” around the MAP value.

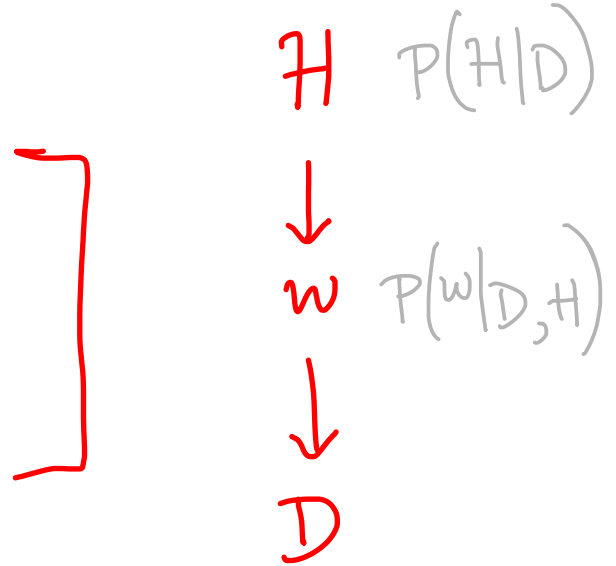


From [MacKay, 2003].

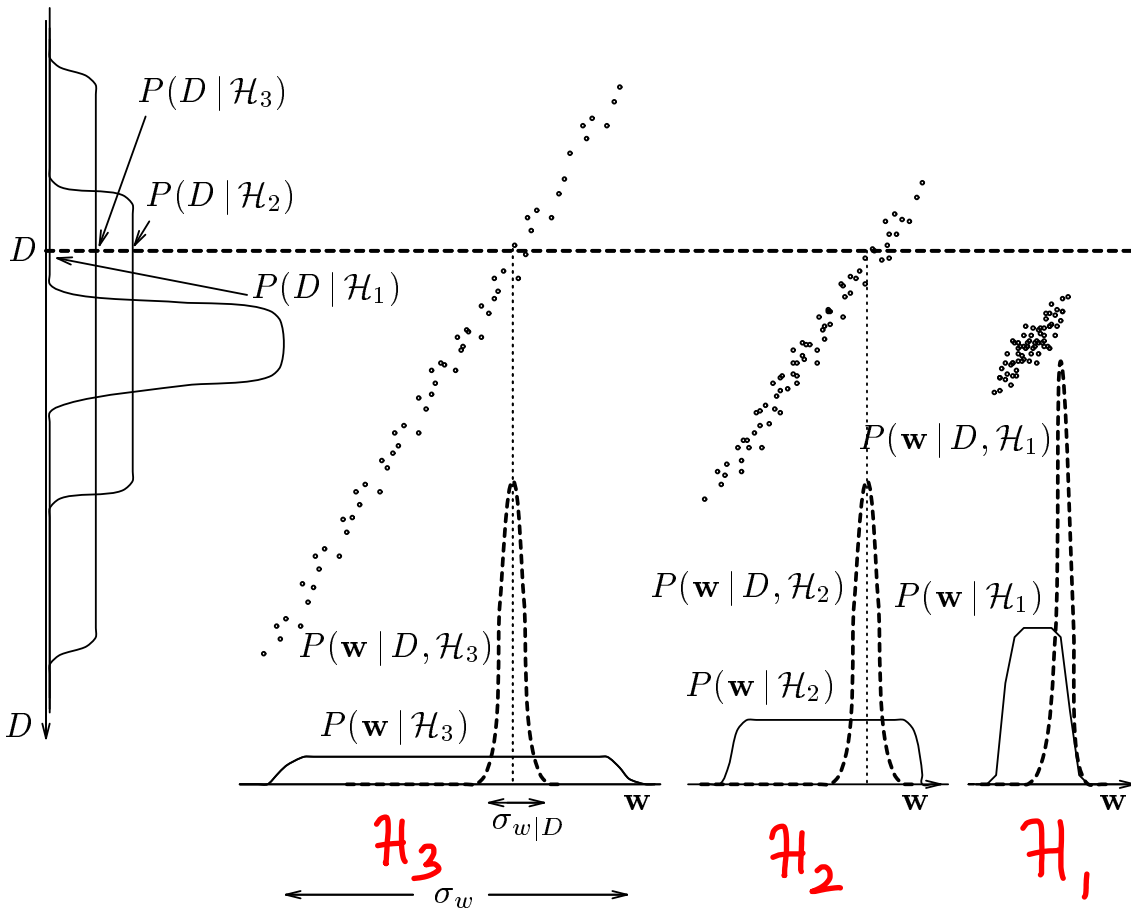
Bayesian inference performs such selection at both model and parameter levels:

$$p(\mathbf{w} | \mathcal{D}, \mathcal{H}_i) = \frac{p(\mathcal{D} | \mathbf{w}, \mathcal{H}_i) p(\mathbf{w} | \mathcal{H}_i)}{p(\mathcal{D} | \mathcal{H}_i)},$$

$$p(\mathcal{H}_i | \mathcal{D}) = \frac{p(\mathcal{D} | \mathcal{H}_i) p(\mathcal{H}_i)}{p(\mathcal{D})}.$$



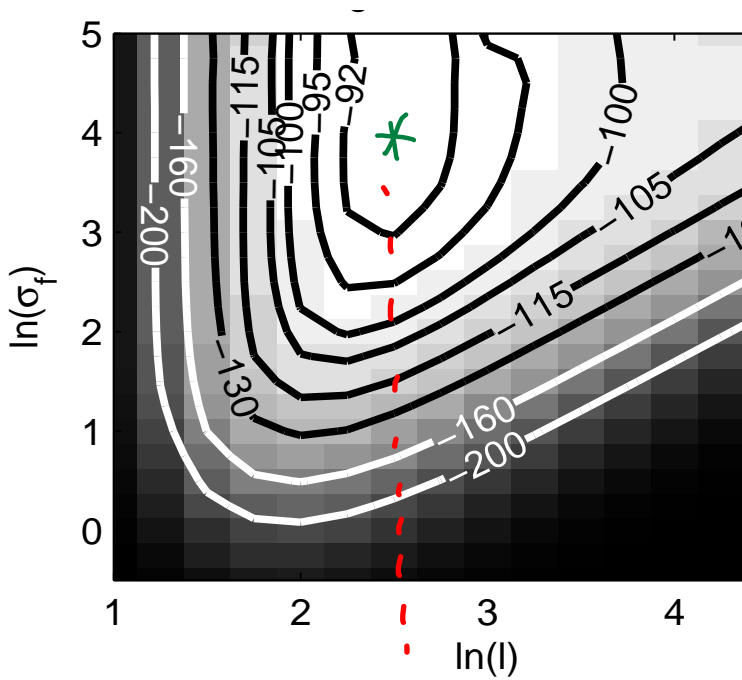
This beautiful illustration by [MacKay, 2003] demonstrates this point for a simple model $\mathcal{D} = w + \text{noise}$. Complexity of the model varies as $\mathcal{H}_3 > \mathcal{H}_2 > \mathcal{H}_1$, but all models have equal prior probabilities.



From [MacKay, 2003].

Demonstration: GP Classification

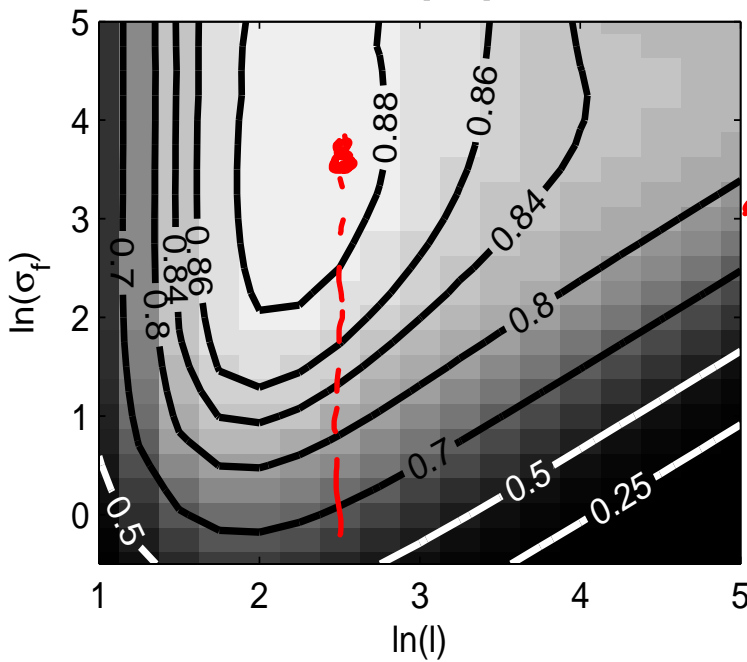
Marginal likelihood on the training data (left) reflects the shape of generalization error evaluated on test data (right).



$$\frac{\|y - Xw\|_2^2}{2} + \lambda \|w\|_2^2$$

$$P(D|\lambda) = \int P(D, w|\lambda) dw$$

→ Marginal likelihood train data

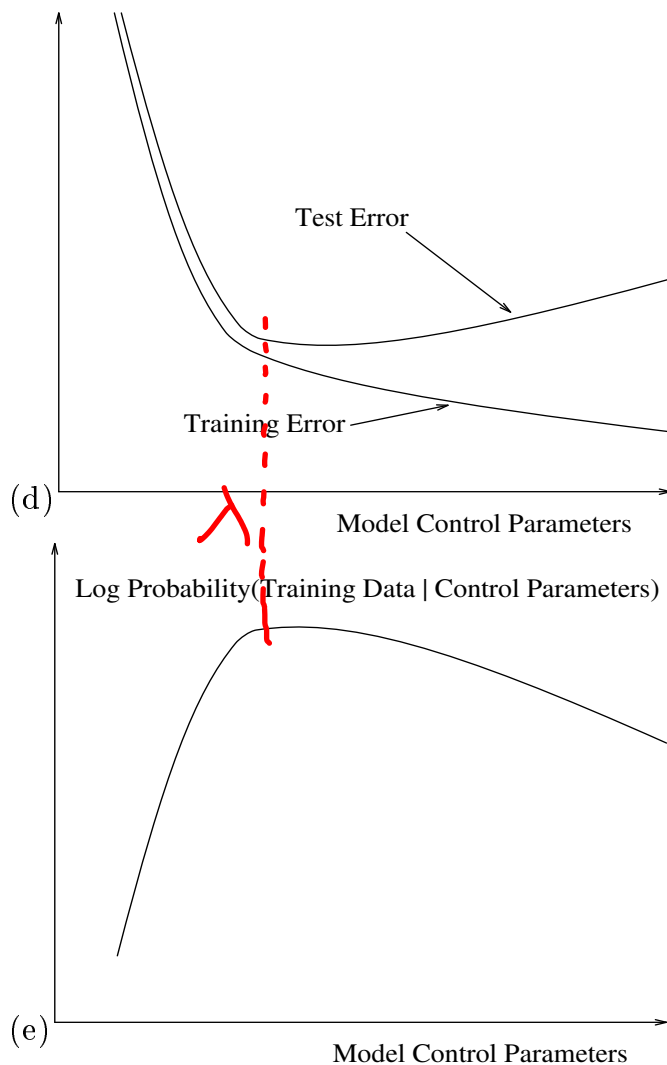
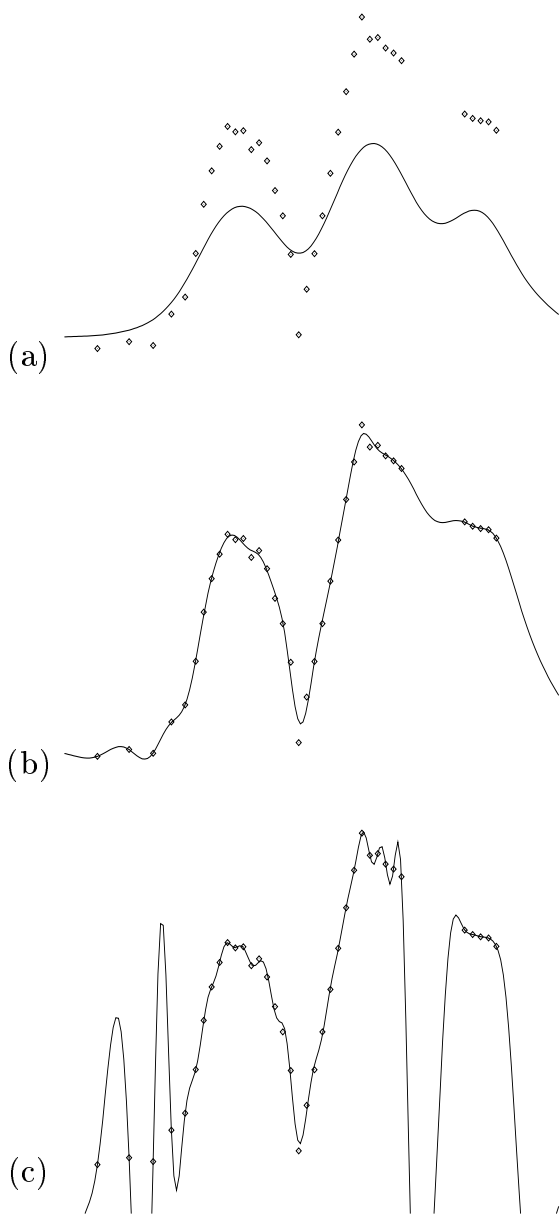


loss on Test data

From [Rasmussen and Williams, 2006].

Demonstration: Neural Networks

When increasing the model complexity, the neural network easily overfits on the test data, but the marginal likelihood on training data reflects the generalization error.



From [MacKay, 2003].

11 Summary

The main computational issue is that we need to average over all possible \mathbf{w} , which is very difficult to do exactly.

$$p(\mathcal{D}) = \int \left[\prod_{n=1}^N p(y_n | f_{\mathbf{w}}(\mathbf{x}_n)) \right] p(\mathbf{w}) d\mathbf{w}.$$

There are at-least three computational challenges:

- 1) Too many factors (large N)
- 2) Too many dimensions (large D)
- 3) Nonconjugacy (e.g., non Gaussian likelihood with a Gaussian prior)

References

- [Bishop, 2006] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*.
- [Khan, 2009] Khan, M. E. (2009). Gaussian likelihood with Gaussian prior on its mean. <https://emtiyaz.github.io/Writings/normNormPdf.pdf> [Online; accessed 21-July-2018].
- [MacKay, 2003] MacKay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge university press.
- [Nickisch and Rasmussen, 2008] Nickisch, H. and Rasmussen, C. (2008). Approximations for binary Gaussian process classification. *Journal of Machine Learning Research*, 9(10).
- [Rasmussen and Williams, 2006] Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press.

List of concepts

gaussian process, 16
gaussian-process prior, 16
regression/classification, 2
conjugate, 14
data-size, 2
data, 2
dimensionality, 2
exponential-family distribution, 3
informative prior, 16
inputs, 2
joint distribution, 5
likelihood, 3
marginal likelihood, 5
maximum a posteriori, 5
maximum-likelihood, 3
output, 2
posterior distribution, 5
prediction, 2
prior distribution, 5
square-exponential kernel, 17

(Notes)

(Notes)