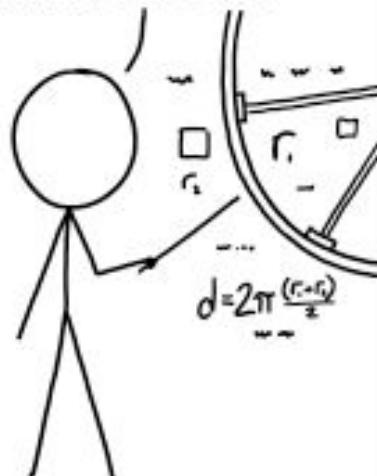


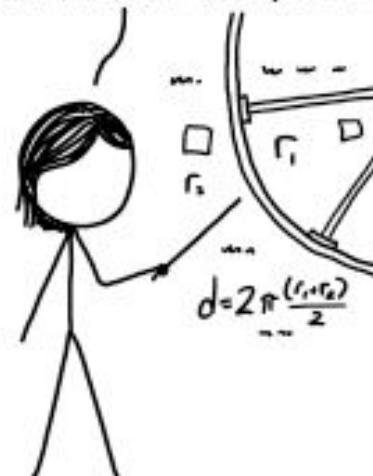
PHYSICIST  
APPROXIMATIONS

WE'LL ASSUME THE  
CURVE OF THIS RAIL  
IS A CIRCULAR ARC  
WITH RADIUS  $R$ .



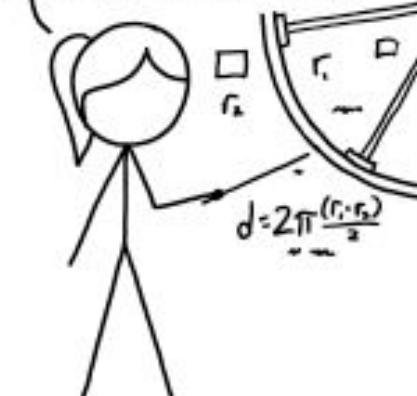
ENGINEER  
APPROXIMATIONS

LET'S ASSUME THIS  
CURVE DEVIATES FROM  
A CIRCLE BY NO MORE  
THAN 1 PART IN 1,000.



COSMOLOGIST  
APPROXIMATIONS

ASSUME PI IS ONE.  
PRETTY SURE IT'S  
BIGGER THAN THAT.  
OK, WE CAN MAKE  
IT TEN. WHATEVER.



# Announcements

Released: Quiz 2, Assignment 2, video assignment

Next three lectures: graphical models

Week 10 Wed: guest lecture

# Approximate Inference + GP Classification

Laplace approximation - in general

Bishop, Chap 4.4, 4.5  
6.4 (6.4.5, part of 6.4.6, 6.4.7)

Laplace approximation - Bayesian logistic regression

GP book *chap 3*  
<http://gaussianprocess.org/gpml/chapters/>

GP classification

Laplace approximation - GP classification

Connection to neural networks

## Laplace approximation in general

[Bishop 4.4]

Goal: find a Gaussian approximation to a probability density defined over a set of continuous variables.

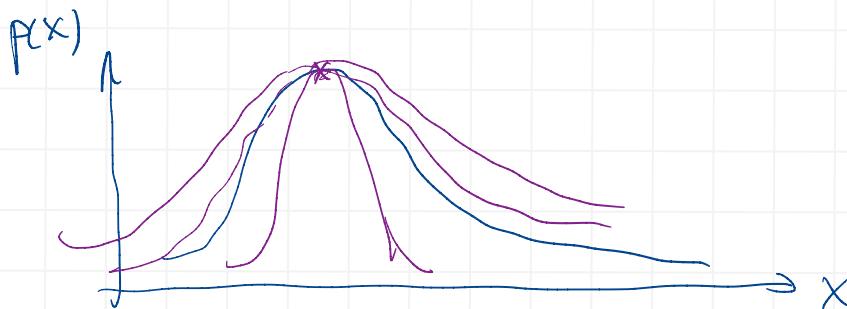
How: find Gaussian pdf  $q(z)$ , centred on a mode of the distribution  $p(z)$

Consider pdf

$$p(z) = \frac{1}{Z} f(z)$$

(4.125)

where  $Z = \int f(z) dz$



# Gaussian pdf $\Leftrightarrow$ log(q) quadratic

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \quad (2.42)$$

$$p(z) = \frac{1}{Z} f(z) \quad (4.125) \quad \text{where } Z = \int f(z) dz$$

Find mode  $z_0$

$$\frac{df(z)}{dz}\Big|_{z=z_0} = 0. \quad (4.126)$$

$$\text{SCM at } z_0: f(z_0) + \frac{1}{1!} f'(z_0) + \frac{1}{2!} f''(z_0)(z - z_0)^2$$

Taylor expansion  
of  $\ln f(z)$  at  $z_0$

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A(z - z_0)^2 \quad (4.127)$$

$$A = -\frac{d^2}{dz^2} \ln f(z)\Big|_{z=z_0}. \quad (4.128)$$

Take exp()

(4.127)

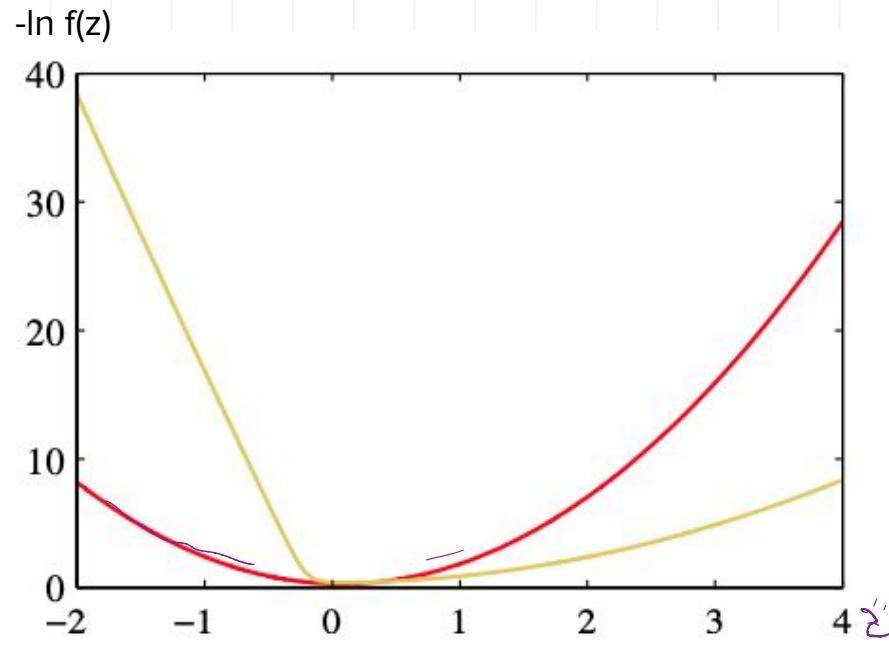
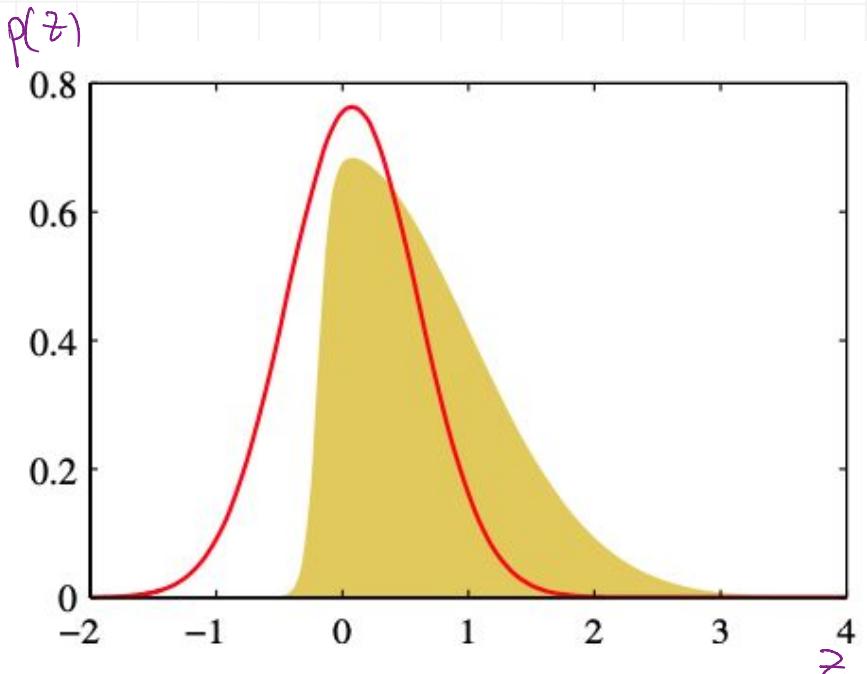
$$f(z) \simeq f(z_0) \exp\left\{-\frac{A}{2}(z - z_0)^2\right\}. \quad (4.129)$$

Normalise to  
obtain  $q(z)$

$$q(z) = \left(\frac{A}{2\pi}\right)^{1/2} \exp\left\{-\frac{A}{2}(z - z_0)^2\right\}. \quad (4.130)$$

Assume  $A > 0$

$$\frac{\partial \ln(f(z))}{\partial z} = \frac{1}{f(z)} \frac{\partial f(z)}{\partial z}$$



**Figure 4.14** Illustration of the Laplace approximation applied to the distribution  $p(z) \propto \exp(-z^2/2)\sigma(20z+4)$  where  $\sigma(z)$  is the logistic sigmoid function defined by  $\sigma(z) = (1 + e^{-z})^{-1}$ . The left plot shows the normalized distribution  $p(z)$  in yellow, together with the Laplace approximation centred on the mode  $z_0$  of  $p(z)$  in red. The right plot shows the negative logarithms of the corresponding curves.

# Laplace approximation in higher dimensions

$$p(\vec{z}) = \frac{1}{Z} f(\vec{z})$$

$f(z)$ , with  $z_0$  a stationary point  $\nabla f(z)|_{z=z_0} = 0$

Taylor expansion of  $\ln f(z)$  around  $z_0$

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2}(z - z_0)^T \mathbf{A}(z - z_0) \quad (4.131)$$

$\xrightarrow{\hspace{1cm}}$   $\begin{matrix} \uparrow \\ \text{matrix} \end{matrix}$   $\boxed{\phantom{0}}$

where the  $M \times M$  Hessian matrix  $\mathbf{A}$  is defined by

$$\mathbf{A} = -\nabla \nabla \ln f(z)|_{z=z_0} \quad (4.132)$$

and  $\nabla$  is the gradient operator. Taking the exponential of both sides we obtain

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{1}{2}(z - z_0)^T \mathbf{A}(z - z_0) \right\}. \quad (4.133)$$

$\uparrow$   
*normalise*

# Laplace approximation in higher dimensions

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\}. \quad (4.133)$$

normalise, use (2.43)

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$$q(\mathbf{z}) = \underbrace{\frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}}}_{\text{normalising constant}} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1}) \quad (4.134)$$

$q(\mathbf{z})$  is a valid multivariate Gaussian distribution iff A positive semi-definite  
→  $\mathbf{z}_0$  is a local maximum, not local min or saddle point

What about  $f(z)$  with multiple modes? Different Laplace approximations for each mode

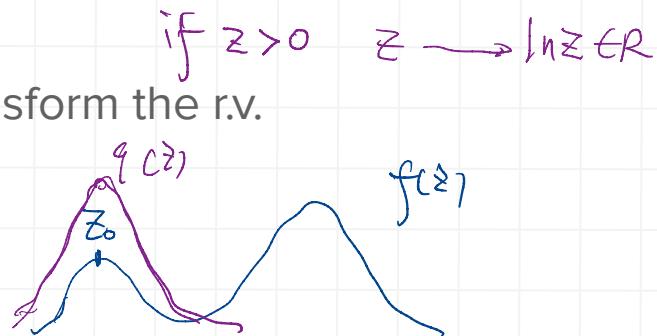
# Pros + cons of Laplace approximation

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z} | \mathbf{z}_0, \mathbf{A}^{-1}) \quad (4.134)$$

- :) Normalisation constant  $Z$  for  $f(z)$  does not need to be known.
- :) CLT → posterior increasingly better approximated by Gaussian as the number of observed data points grow.

:(  
Assumes the domain of  $z$  is  $\mathbb{R}$ , or  $\mathbb{R}^d$ , may need to transform the r.v.

:(  
Is based purely on  $f(z)$  around  $z_0$ , no global info



# Approximate Inference + GP Classification

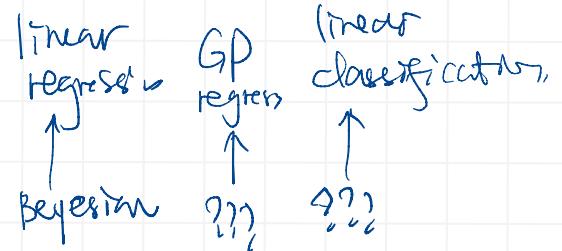
Laplace approximation - in general

Laplace approximation - Bayesian logistic regression

GP classification

Laplace approximation - GP classification

Connection to neural networks





## Recap: Logistic regression

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

data set  $\{\phi_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $\phi_n = \phi(\mathbf{x}_n)$ .

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad (4.89)$$

$\sigma(\omega^T \phi)$  OR  $1 - \sigma(\omega^T \phi)$

Negative log-likelihood, or cross-entropy error function

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4.90)$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n \quad (4.91)$$

# Bayesian logistic regression

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} \quad (4.90)$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \quad (4.140)$$

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{w}) \quad (4.141) \quad \xleftarrow{\text{PCT}} \text{const}$$

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{t}) &= -\frac{1}{2} \underbrace{(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)}_{\text{quadratic part of } \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)} \\ &\quad + \sum_{n=1}^N \underbrace{\{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}}_{\text{come from } E(\mathbf{w})} + \text{const} \quad (4.142) \end{aligned}$$

normaliser of  $\mathcal{N}(\dots)$

# Laplace approximation for Bayesian logistic regression

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_{n=1}^N \underbrace{\{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}}_{\textcircled{a}} + \text{const} \quad (4.142) \end{aligned}$$

$$y(\phi) = \sigma(\mathbf{w}^T \phi)$$

① compute local max  $\mathbf{z}_0$

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N). \quad (4.144)$$

② compute Hessian around  $\mathbf{z}_0$

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N \underbrace{y_n(1-y_n)}_{\textcircled{b}} \underbrace{(\mathbf{\phi}_n \mathbf{\phi}_n^T)}_{\textcircled{c}}. \quad (4.143)$$

$$y(\phi) = \sigma(\mathbf{w}^T \phi)$$

$$\rightarrow \frac{d\sigma}{d\phi} = \sigma(1-\sigma)$$

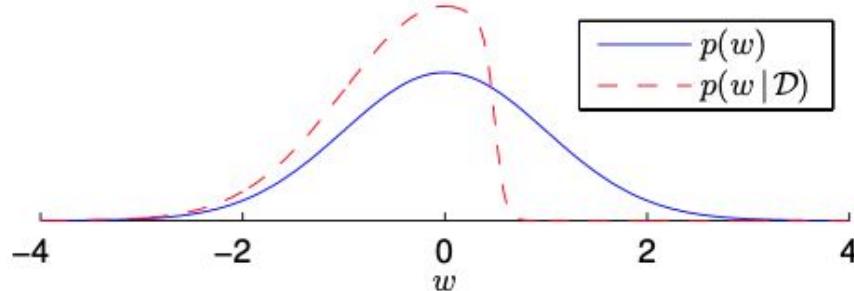
$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \mathbf{\phi}_n$$

$$\begin{aligned} \textcircled{b} \quad & \nabla_{\mathbf{w}} t_n \ln \sigma(\mathbf{w}^T \phi_n) \\ &= t_n \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \frac{1}{\sigma(\mathbf{w}^T \phi_n)} \phi_n \end{aligned}$$

## Example 1

$$p(w) \propto \mathcal{N}(w; 0, 1)$$

$$p(w | \mathcal{D}) \propto \mathcal{N}(w; 0, 1) \sigma(10 - 20w).$$

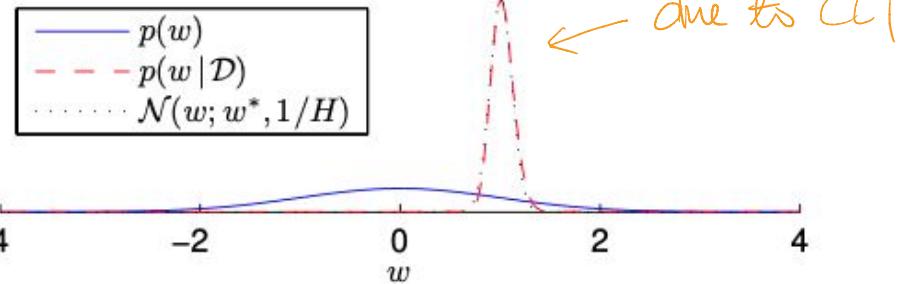


## Example 2

As another example, I generated  $N = 500$  labels,  $\{z^{(n)}\}$ , from a logistic regression model with no bias and with  $w=1$  at  $x^{(n)} \sim \mathcal{N}(0, 10^2)$ . Then,

$$p(w) \propto \mathcal{N}(w; 0, 1)$$

$$p(w | \mathcal{D}) \propto \mathcal{N}(w; 0, 1) \prod_{n=1}^{500} \sigma(wx^{(n)}z^{(n)}), \quad z^{(n)} \in \{\pm 1\}.$$

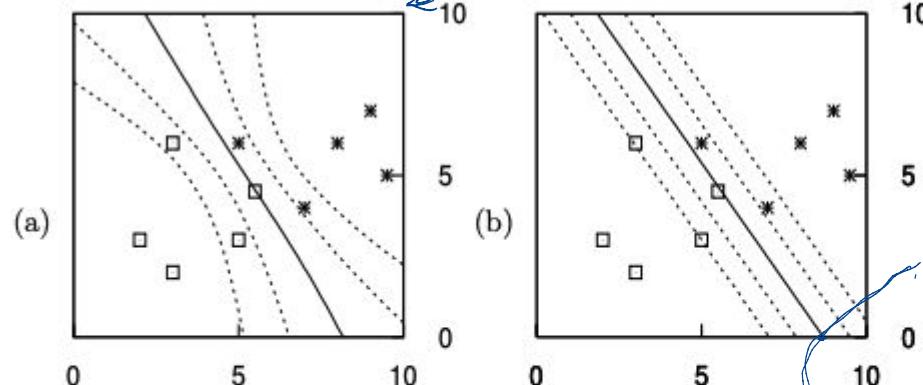


# What does the predictive distribution look like?

$$p(\mathcal{C}_1|\phi, \mathbf{t}) = \int p(\mathcal{C}_1|\phi, \mathbf{w}) p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad (4.145)$$

laplace approx.

$$= \int \sigma(a) \mathcal{N}(a|\mu_a, \sigma_a^2) da.$$



Again, the axes are the input features  $x_1$  and  $x_2$ . The right hand figure shows  $P(y=1 | \mathbf{x}, \mathbf{w}^*)$  for some fitted weights  $\mathbf{w}^*$ . No matter how these fitted weights are chosen, the contours have to be linear. The parallel contours mean that the uncertainty of predictions falls at the same rate when moving away from the decision boundary, no matter how far we are from the training inputs.

assumes  $p(\mathbf{w}|\mathbf{t}) = \delta(\mathbf{w}^*)$

$\sigma(\mathbf{w}^T \phi)$

over  $\phi = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

# Approximate Inference + GP Classification

Laplace approximation - in general

Laplace approximation - Bayesian logistic regression

GP classification

Laplace approximation - GP classification

Connection to neural networks

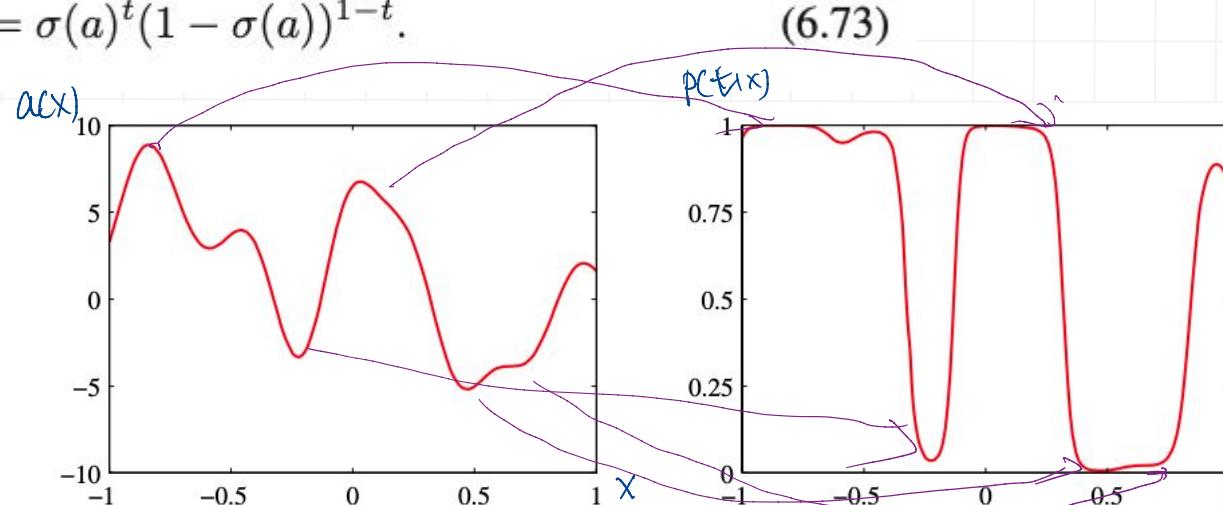
# GP for classification

$$\mathbf{a}(x) \sim GP(\mathbf{0}, K)$$

mean fn  
 $\mathbf{0}$   
covariance fn

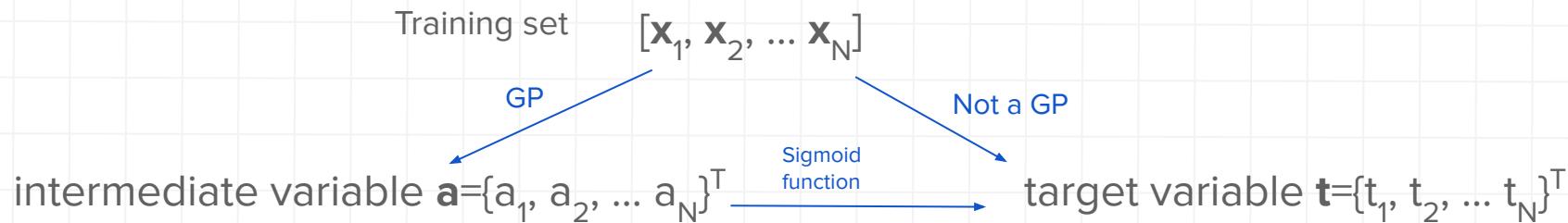
Training set: input  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , target variable  $\mathbf{t} = [t_1, t_2, \dots, t_N]^T$ , with  $t_i \in \{0, 1\}$

$$p(t|a) = \sigma(a)^t(1 - \sigma(a))^{1-t}.$$



**Figure 6.11** The left plot shows a sample from a Gaussian process prior over functions  $a(x)$ , and the right plot shows the result of transforming this sample using a logistic sigmoid function.

# What we just said



For regression

$$p(t_{N+1}) = \mathcal{N}(t_{N+1} | \mathbf{0}, \mathbf{C}_{N+1}) \quad (6.64)$$

Rename the variable

$$p(a_{N+1}) = \mathcal{N}(a_{N+1} | \mathbf{0}, \mathbf{C}_{N+1}). \quad (6.74)$$

## Covariance function of the GP $\mathbf{a}$

$$p(\mathbf{a}_{N+1}) = \mathcal{N}(\mathbf{a}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1}). \quad (6.74)$$

last week

$$C(x_n, x_m) = f(x_n, x_m)$$

$$+ \beta^2 \delta_{nm}.$$

↑  
noise -

- Assume it's noise-less
- BUT add diagonal term to ensure that it's positive semi-definite

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \nu \delta_{nm} \quad (6.75)$$

Assume kernel function  $k(\mathbf{x}, \mathbf{x}' ; \Theta)$

# How to do prediction?

Assume:  $X_{N+1}$  hidden.

Desired:

$$p(t_{N+1} = 1 | \mathbf{t}_N)$$

Have:

$$p(t_{N+1} = 1 | a_{N+1}) = \sigma(a_{N+1})$$

Plug in GP prediction

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (6.66)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \quad (6.67)$$

$x$  is “hidden” inside kernels  $\mathbf{k}$  and  $\mathbf{C}$

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}(a_{N+1} | \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}). \quad (6.78)$$

Use Bayes rule:

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \quad (6.76)$$

*assumes  $t_{N+1} \perp t_n$  given  $a_{N+1}$*

# Computing posterior for GP classification

*predictive*

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \quad (6.76)$$

want to  
use

$p(a_{N+1} | \mathbf{a}_N)$ :

$$\begin{aligned} p(a_{N+1} | \mathbf{t}_N) &= \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \quad \text{→ Bayes rule,} \\ &\xrightarrow{\text{const}} \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N \\ &= \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \end{aligned} \quad (6.77)$$

$(x_n, t_n)$  known  
 $p(\mathbf{a}_N)$  prior known.

Laplace approximation

$$\rightarrow p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}(a_{N+1} | \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}). \quad (6.78)$$

Want to approximate  $\underline{p(\mathbf{a}_N | \mathbf{t}_N)}$

$\rightarrow$  logistic reg.  
 $p(\mathbf{t}_N | \text{land})$

$$\underline{\underline{P(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{0}, \mathbf{C}_N)}}$$

$$p(\mathbf{t}_N | \mathbf{a}_N) = \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n). \quad (6.79)$$

*expand + simplify.*

Taylor expansion of  $\underline{p(\mathbf{a}_N | \mathbf{t}_N)}$

$$\begin{aligned} \Psi(\mathbf{a}_N) &= \underline{\ln p(\mathbf{a}_N)} + \underline{\ln p(\mathbf{t}_N | \mathbf{a}_N)} \\ &= -\frac{1}{2} \underline{\mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N} - \frac{N}{2} \underline{\ln(2\pi)} - \frac{1}{2} \underline{\ln |\mathbf{C}_N|} + \underline{\mathbf{t}_N^T \mathbf{a}_N} \\ &\quad - \sum_{n=1}^N \underline{\ln(1 + e^{a_n})} + \underline{\text{const.}} \end{aligned} \quad (6.80)$$

const  
 $-\ln p(\mathbf{t}_N)$

$(1 + e^{-a_n})$

# Laplace approximation of $p(\mathbf{a}_N | \mathbf{t}_N)$

$$\begin{aligned}\Psi(\mathbf{a}_N) &= \ln p(\mathbf{a}_N) + \ln p(\mathbf{t}_N | \mathbf{a}_N) \\ &= -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N \\ &\quad - \sum_{n=1}^N \ln(1 + e^{a_n}) + \text{const.}\end{aligned}\tag{6.80}$$

$$\nabla \Psi(\mathbf{a}_N) = \mathbf{t}_N - \sigma_N - \mathbf{C}_N^{-1} \mathbf{a}_N \tag{6.81}$$

$$\nabla \nabla \Psi(\mathbf{a}_N) = -\mathbf{W}_N - \mathbf{C}_N^{-1} \tag{6.82}$$

Claim:  $-\nabla \nabla \Psi(\mathbf{a}_N)$  positive definite, but posterior is non-Gaussian (since  $\mathbf{W}_N$  depend on  $\mathbf{a}_N$ )

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1}). \tag{6.86}$$

$\downarrow$   
valid Laplace approximation,

$$x_1, \dots, x_n \rightarrow a_1, \dots, a_m \rightarrow b_1, \dots, b_n$$

(a) find local max scalar

$$(b) H = -\nabla \nabla$$

$$\nabla \Psi(\mathbf{a}_N) = 0 \Rightarrow \mathbf{a}_N^* = \mathbf{t}_N - \sigma_N$$

$$\sigma_N = [\sigma(a_n)] \quad \frac{\partial \sigma}{\partial a} = f(1-f)$$

$$\mathbf{W}_N = \underbrace{\text{diag}}_{\text{symmetric}}[\underbrace{\sigma(a_n)}_{(0, 1)}, \underbrace{1 - \sigma(a_n)}_{(0, 1)}]$$

$$\mathbf{a}_N^* = \mathbf{C}_N (\mathbf{t}_N - \sigma_N). \tag{6.84}$$

$$\mathbf{H} = -\nabla \nabla \Psi(\mathbf{a}_N) = \mathbf{W}_N + \mathbf{C}_N^{-1} \tag{6.85}$$

Q: What do we use for  $f(a_n)$

# Bringing it back together

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \quad (6.76)$$

$$\int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \quad (6.77)$$

$$\sigma(a_{N+1})$$

$$\boxed{p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T)} \quad (2.115)$$

$$\mathcal{N}(a_{N+1} | \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k})$$

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1}).$$

$$\mathbf{a}_N^* = \mathbf{C}_N(\mathbf{t}_N - \boldsymbol{\sigma}_N).$$

$$\mathbf{H} = -\nabla \nabla \Psi(\mathbf{a}_N) = \mathbf{W}_N + \mathbf{C}_N^{-1}$$

$$\underline{\mathbb{E}[a_{N+1} | \mathbf{t}_N]} = \mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N) \quad (6.87)$$

$$\underline{\text{var}[a_{N+1} | \mathbf{t}_N]} = c - \mathbf{k}^T (\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1} \mathbf{k}. \quad (6.88)$$

## Bringing it back together

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \quad (6.76)$$

$$\sigma(a_{N+1})$$

$$\mathbb{E}[a_{N+1} | \mathbf{t}_N] = \mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N) \quad (6.87)$$

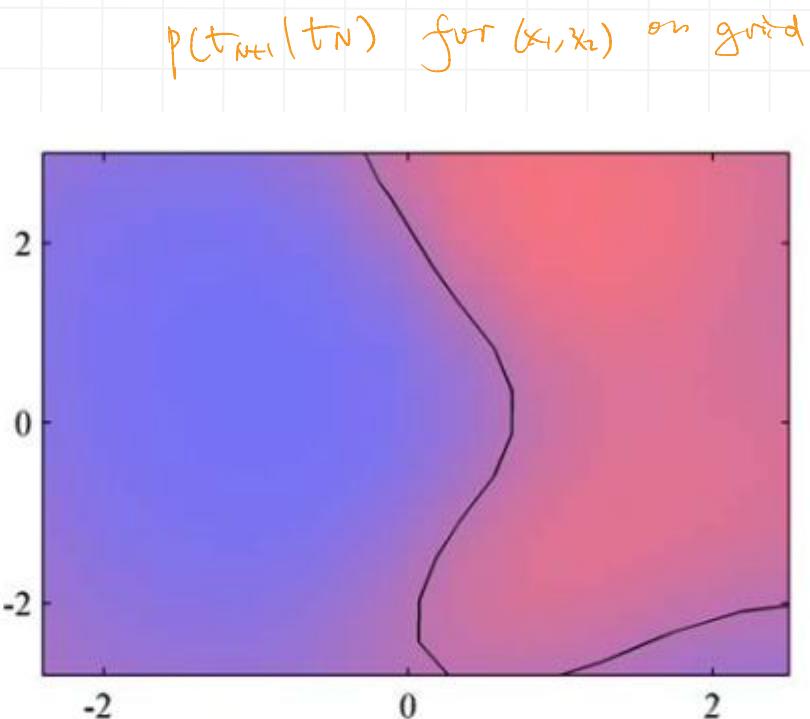
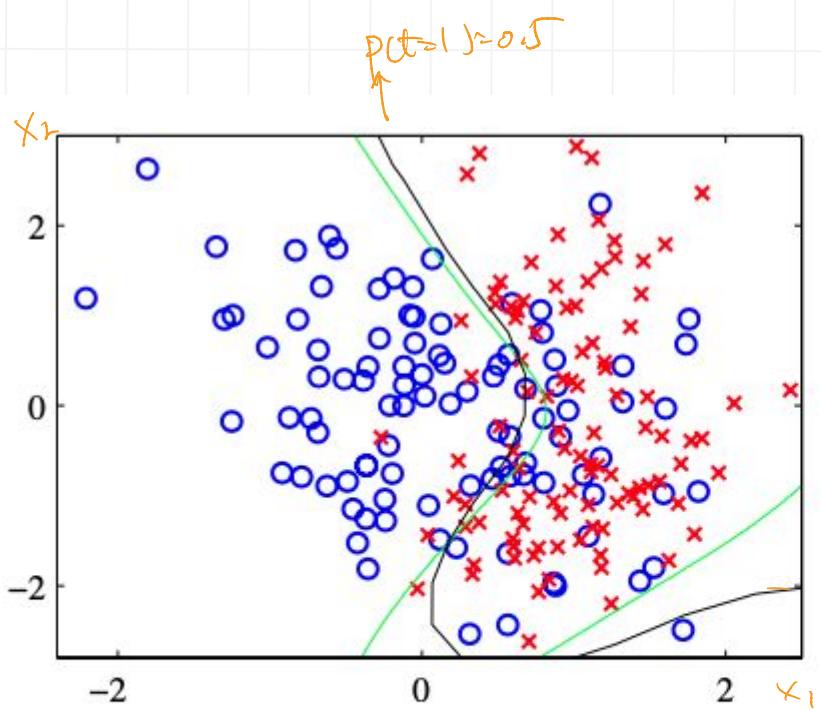
$$\text{var}[a_{N+1} | \mathbf{t}_N] = c - \mathbf{k}^T (\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1} \mathbf{k}. \quad (6.88)$$

$$\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da \simeq \sigma(\kappa(\sigma^2)\mu) \quad (4.153)$$

$$\kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2}.$$

$$c - \mathbf{k}^T (\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1} \mathbf{k}$$

$$\mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N)$$



**Figure 6.12** Illustration of the use of a Gaussian process for classification, showing the data on the left together with the optimal decision boundary from the true distribution in green, and the decision boundary from the Gaussian process classifier in black. On the right is the predicted posterior probability for the blue and red classes together with the Gaussian process decision boundary.

# Connection to neural nets

2-layer neural network, with  $M$  hidden units

Bayesian neural network function  $f(x, \mathbf{w})$ , with prior over  $\mathbf{w}$

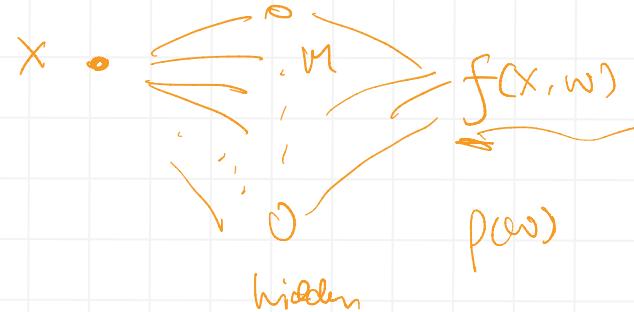
(Neal 1996) for a broad class of prior distributions over  $\mathbf{w}$ , the distribution of functions generated by a neural network will tend to a Gaussian process when  $M$  tends to infinity.

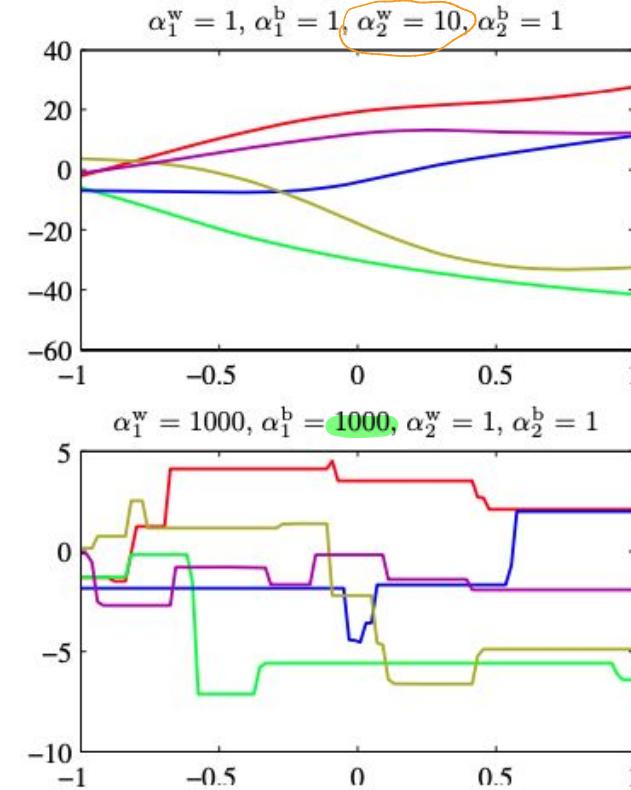
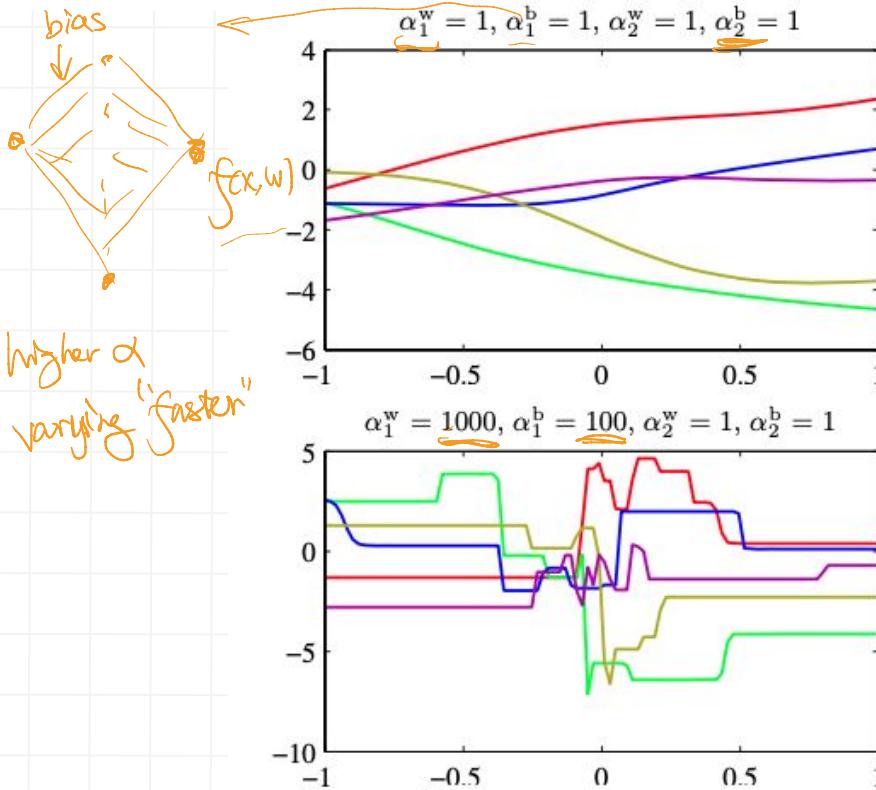
$$\text{e.g. } k(x' - x) = f(\|x - x'\|) = \exp(-\frac{\sigma}{2} \|x - x'\|^2)$$

*non-stationary*  $k(x, x') = x^T x'$

Covariance function  $k(x, x')$  non-stationary for neural nets with probit and Gaussian activations.

Weight prior determine the lengths scales of the neural net function.





**Figure 5.11** Illustration of the effect of the hyperparameters governing the prior distribution over weights and biases in a two-layer network having a single input, a single linear output, and 12 hidden units having ‘tanh’ activation functions. The priors are governed by four hyperparameters  $\alpha_1^b$ ,  $\alpha_1^w$ ,  $\alpha_2^b$ , and  $\alpha_2^w$ , which represent the precisions of the Gaussian distributions of the first-layer biases, first-layer weights, second-layer biases, and second-layer weights, respectively. We see that the parameter  $\alpha_2^w$  governs the vertical scale of functions (note the different vertical axis ranges on the top two diagrams),  $\alpha_1^w$  governs the horizontal scale of variations in the function values, and  $\alpha_1^b$  governs the horizontal range over which variations occur. The parameter  $\alpha_2^b$ , whose effect is not illustrated here, governs the range of vertical offsets of the functions.

# Gaussian Processes 2

GP classification

Laplace approximation - in general

Laplace approximation - Bayesian logistic regression

Laplace approximation - GP classification