

# PROBABILISTIC INFERENCE AND LEARNING

## LECTURE 10

### GAUSSIAN PROCESS CLASSIFICATION AND OTHER NON-GAUSSIAN OBSERVATION MODELS

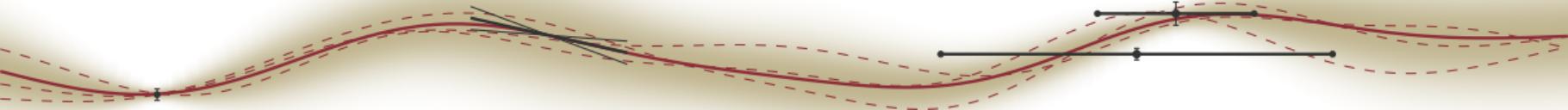
Philipp Hennig

19 November 2018

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN

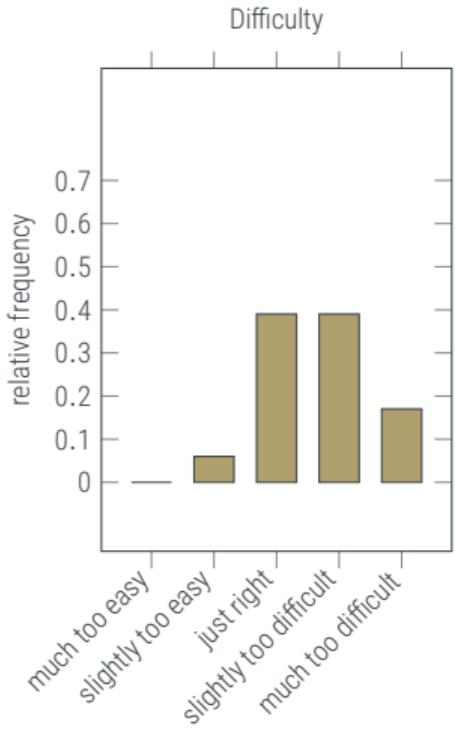
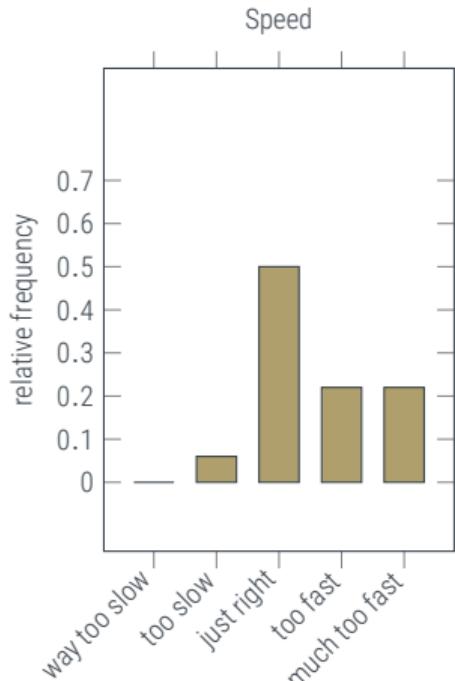
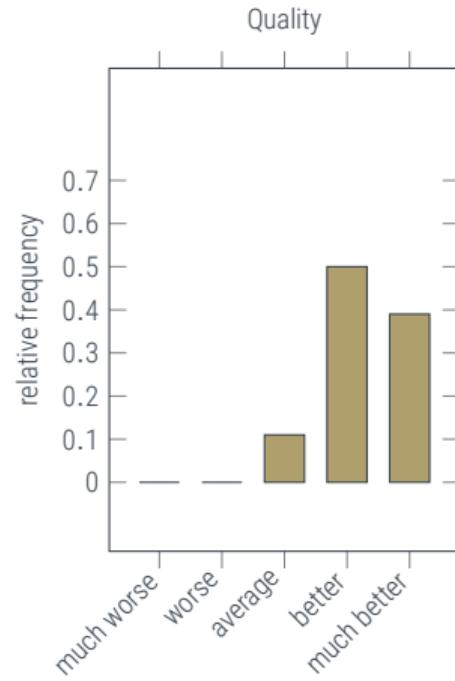


FACULTY OF SCIENCE  
DEPARTMENT OF COMPUTER SCIENCE  
CHAIR FOR THE METHODS OF MACHINE LEARNING



# Last Lecture: Debrief

Feedback dashboard





# Last Lecture: Debrief

## Detailed Feedback

### Things you did not like:

- ❖ “please no break [on Wednesday]”
- ❖ Lecture has to end by 10:50 at latest
- ❖ “can you just report the results and let us do the derivations?”
- ❖ long derivation
- ❖ lecture hall is too cold
- ❖ the code should contain equations!
- ❖ use something other than color to highlight

### Things you did not understand:

- ❖ “I’m not sure what we’re assuming for the Kalman filter to work”
- ❖ why
$$\int p(x_t | x_{t+1}, Y_{0:t}) dx_t = 1?$$

### Things you enjoyed:

- ❖ detailed derivations
- ❖ plots
- ❖ code
- ❖ slides
- ❖ “the concepts – really getting into it”
- ❖ recaps and connections (but too long now)



## Overview of Lectures so far:

0. Introduction to Reasoning under Uncertainty
1. Probabilistic Reasoning
2. Probabilities over Continuous Variables
3. Gaussian Probability Distributions
4. Gaussian Parametric Regression
5. More on Parametric Regression
6. Gaussian Processes
7. More on Kernels & GPs
8. A practical GP example
9. Markov Chains, Time Series, Filtering
10. Classification (TODAY)
11. Empirical Example of Classification (WED)
12. More Connections to SLT (NEXT MON)
13. Theory of Filtering / SDEs (NEXT WED)

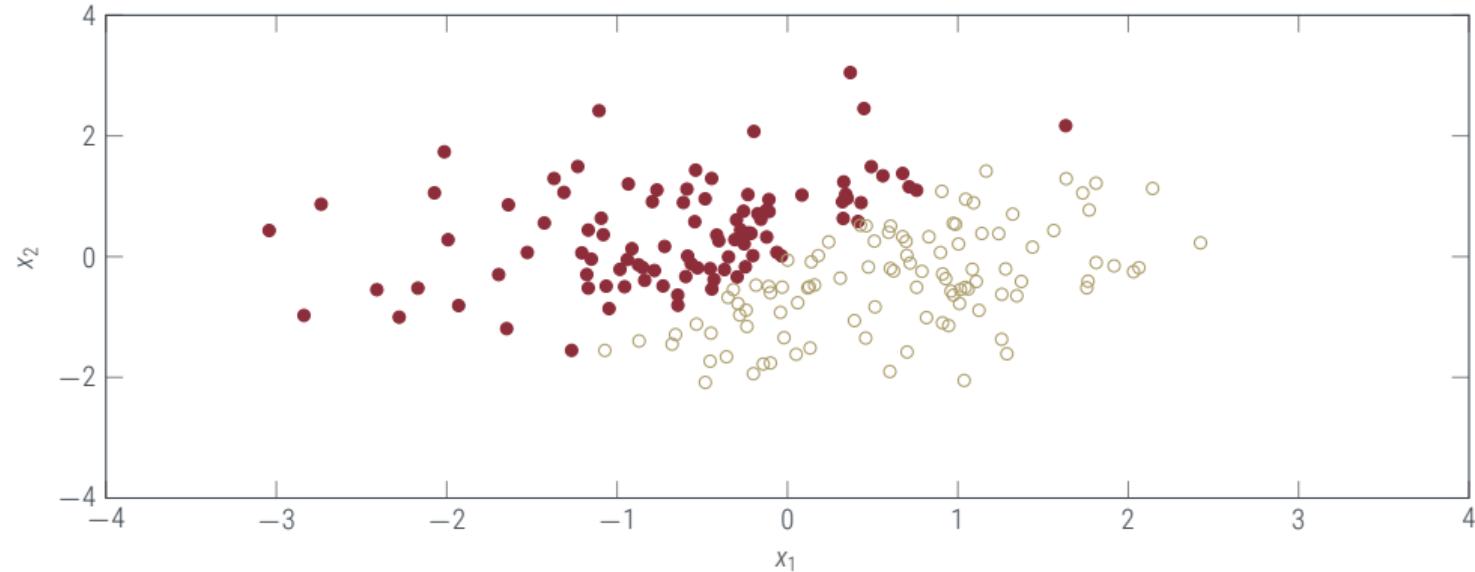
## Today:

- Classification: Supervised Learning with Discrete Outputs
- *Conceptually*, a simple generalization of GP regression
- *Computationally*, we will have to introduce new approximation methods



# Classification Problems

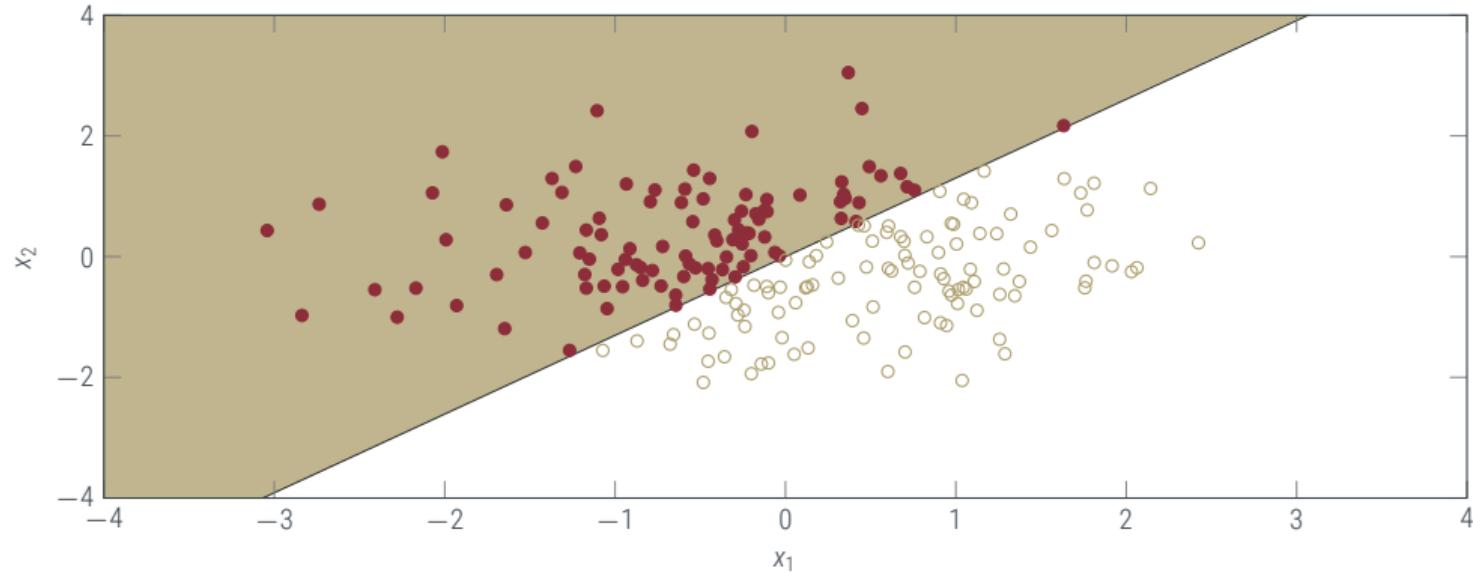
a typography of supervised learning problems





# Classification Problems

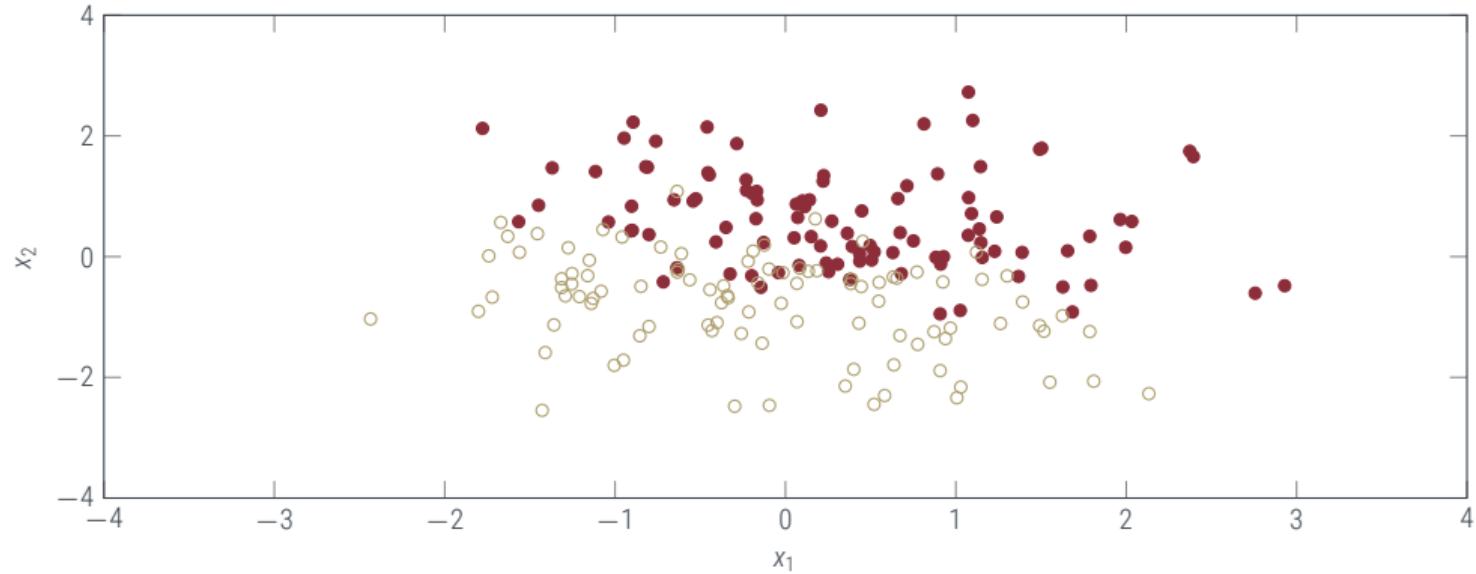
a typography of supervised learning problems





# Classification Problems

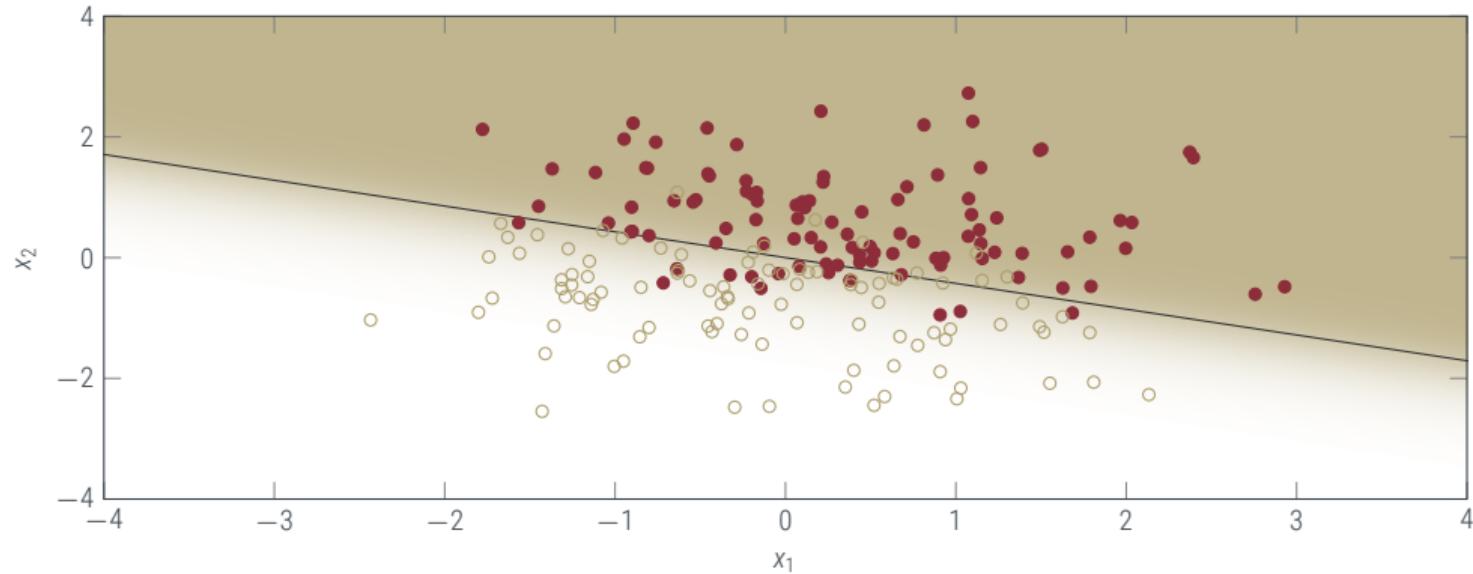
a typography of supervised learning problems





# Classification Problems

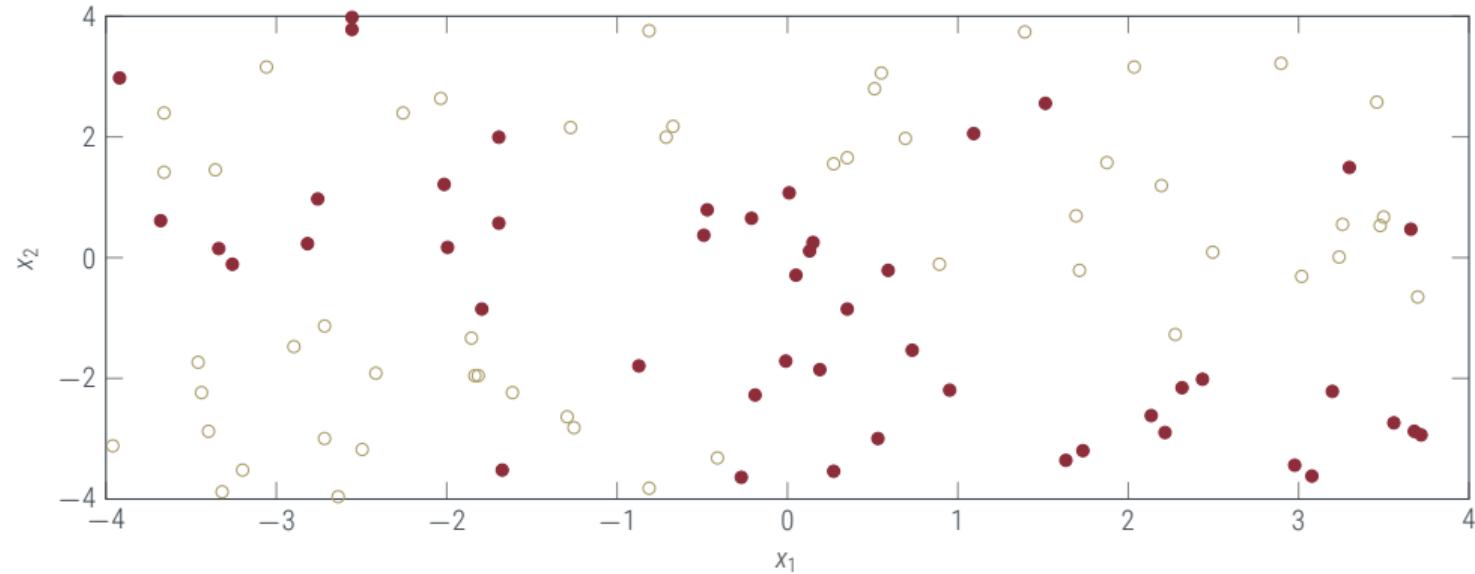
a typography of supervised learning problems





# Classification Problems

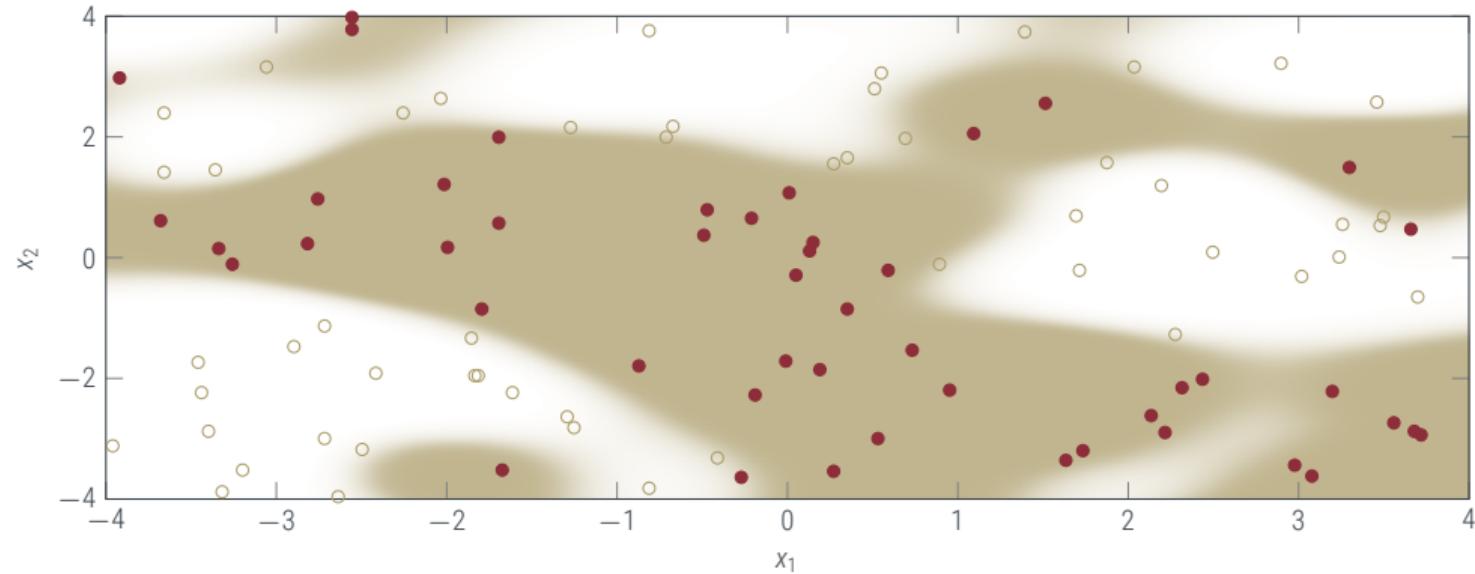
a typography of supervised learning problems





# Classification Problems

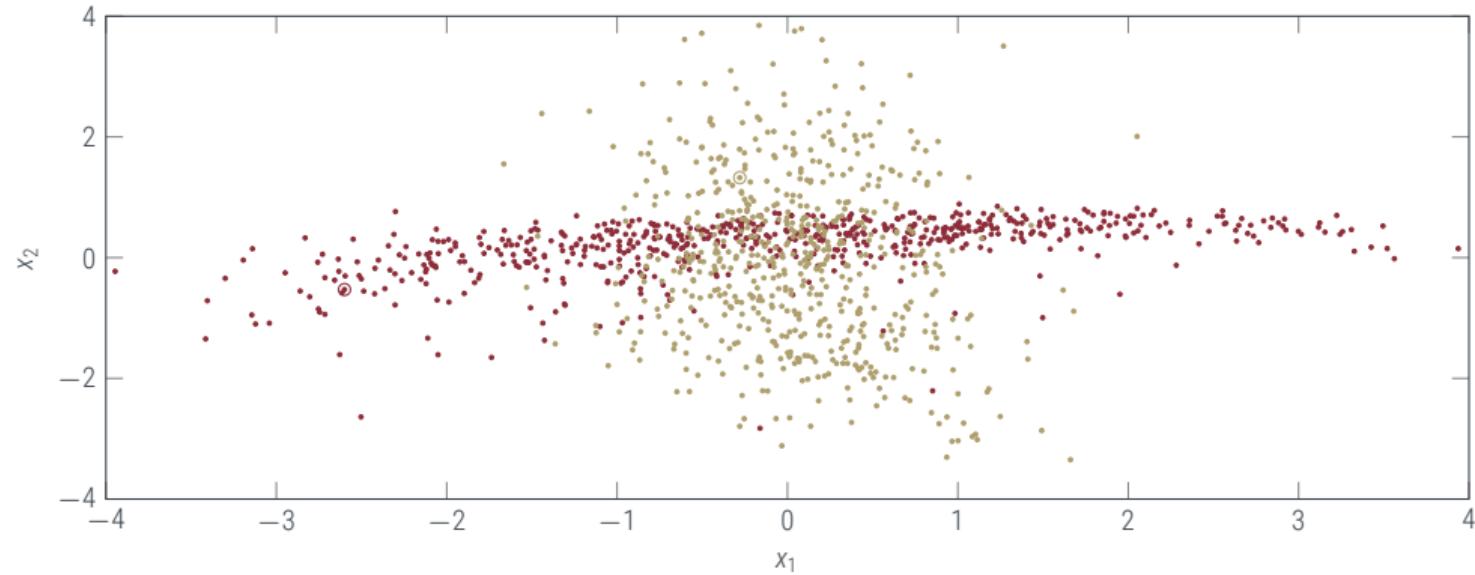
a typography of supervised learning problems





# Classification Problems

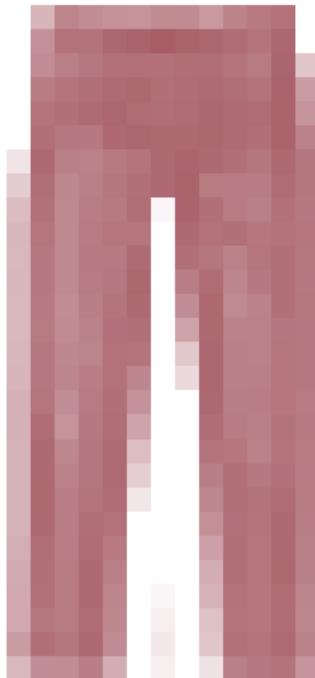
a typography of supervised learning problems





# Classification Problems

a typography of supervised learning problems



<https://github.com/zalandoresearch/fashion-mnist>



# Classification vs. Regression

Two types of supervised learning problems

## Regression:

Given supervised *data* (special case  $d = 1$ : univariate regression)

$$(X, Y) := (x_i, y_i)_{i=1, \dots, n} \text{ with } x_i \in \mathbb{X}, y_i \in \mathbb{R}^d$$

find function  $f : \mathbb{X} \rightarrow \mathbb{R}^d$  such that  $f$  "models"  $Y \approx f(X)$ .

## Classification:

Given supervised *data* (special case  $m = 1$ : binary classification)

$$(X, Y) := (x_i, c_i)_{i=1, \dots, n} \text{ with } x_i \in \mathbb{X}, c_i \in \{0, \dots, d - 1\}$$

find probability  $\pi : \mathbb{X} \rightarrow U^d$  ( $U^d = \{p \in [0, 1]^d : \sum_{i=1}^d p_i = 1\}$ ) such that  $\pi$  "models"  $y_i \sim \pi_{x_i}$ .

Classification and regression are quite similar at first sight. But classification is a bit more subtle...



Until further notice, consider only discriminative **binary** classification:

$$y \in \{-1; +1\} \quad x \rightarrow \pi(x) =: \pi_x \in [0, 1]$$

$$p(y | x) = \begin{cases} \pi(x) & \text{if } y = 1 \\ 1 - \pi(x) & \text{if } y = -1 \end{cases}$$



Until further notice, consider only discriminative **binary** classification:

$$y \in \{-1; +1\} \quad x \rightarrow \pi(x) =: \pi_x \in [0, 1]$$

$$p(y | x) = \begin{cases} \pi(x) & \text{if } y = 1 \\ 1 - \pi(x) & \text{if } y = -1 \end{cases}$$

Discriminative learning phrased probabilistically:

- We would like to *learn*  $\pi_x(y) = p(y | x)$



Until further notice, consider only discriminative **binary** classification:

$$y \in \{-1; +1\} \quad x \rightarrow \pi(x) =: \pi_x \in [0, 1]$$

$$p(y | x) = \begin{cases} \pi(x) & \text{if } y = 1 \\ 1 - \pi(x) & \text{if } y = -1 \end{cases}$$

Discriminative learning phrased probabilistically:

- We would like to *learn*  $\pi_x(y) = p(y | x)$
- This is *almost* like regression:  $p(y | x) = \mathcal{N}(y; f_x, \sigma^2) = \pi_x(y)$

Until further notice, consider only discriminative **binary** classification:

$$y \in \{-1; +1\} \quad x \rightarrow \pi(x) =: \pi_x \in [0, 1]$$

$$p(y | x) = \begin{cases} \pi(x) & \text{if } y = 1 \\ 1 - \pi(x) & \text{if } y = -1 \end{cases}$$

Discriminative learning phrased probabilistically:

- We would like to *learn*  $\pi_x(y) = p(y | x)$
- This is *almost* like regression:  $p(y | x) = \mathcal{N}(y; f_x, \sigma^2) = \pi_x(y)$
- only the *domain* is wrong:  $y \in \{-1; 1\}$  vs.  $y \in \mathbb{R}$ .



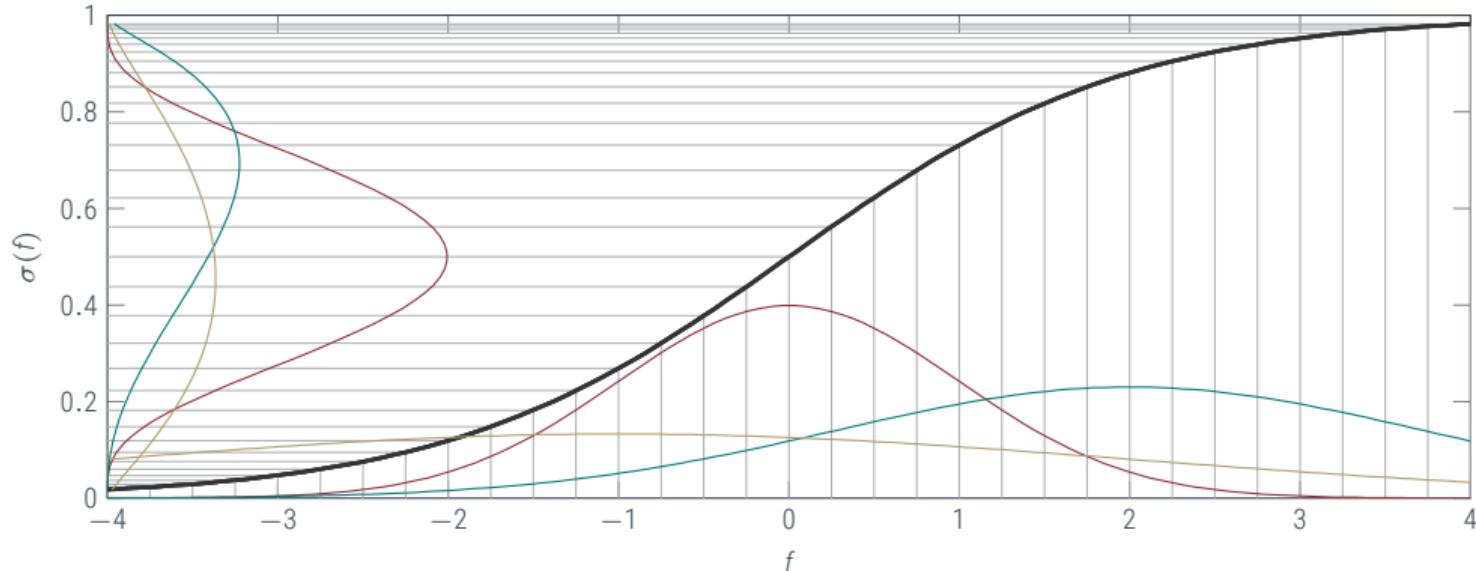
# Let's not throw out Gauss just yet

Turning Gaussian process models into discrete likelihoods



# Link Functions

sigmoids – the logistic link function



$$\pi_f = \sigma(f) = \frac{1}{1 + \exp(-f)} = \int_{-\infty}^f \frac{1}{4} \operatorname{sech}^2\left(\frac{x}{2}\right) dx$$

$$\sigma(f) = 1 - \sigma(-f) \quad f(\pi) = \ln \pi - \ln(1 - \pi) \quad \frac{d\pi}{df} = \pi(f) \cdot (1 - \pi(f))$$



# CODE

Generative Model for Logistic Regression.ipynb



# A Gaussian Process model for Classification

## Logistic Regression

$$p(f) = \mathcal{GP}(f; m, k)$$

$$p(y | f_x) = \sigma(yf_x) = \begin{cases} \sigma(f) & \text{if } y = 1 \\ 1 - \sigma(f) & \text{if } y = -1 \end{cases} \quad \text{using } \sigma(x) = 1 - \sigma(-x).$$

# A Gaussian Process model for Classification

## Logistic Regression

$$p(f) = \mathcal{GP}(f; m, k)$$

$$p(y | f_x) = \sigma(yf_x) = \begin{cases} \sigma(f) & \text{if } y = 1 \\ 1 - \sigma(f) & \text{if } y = -1 \end{cases} \quad \text{using } \sigma(x) = 1 - \sigma(-x).$$

The problem: The posterior is not Gaussian!

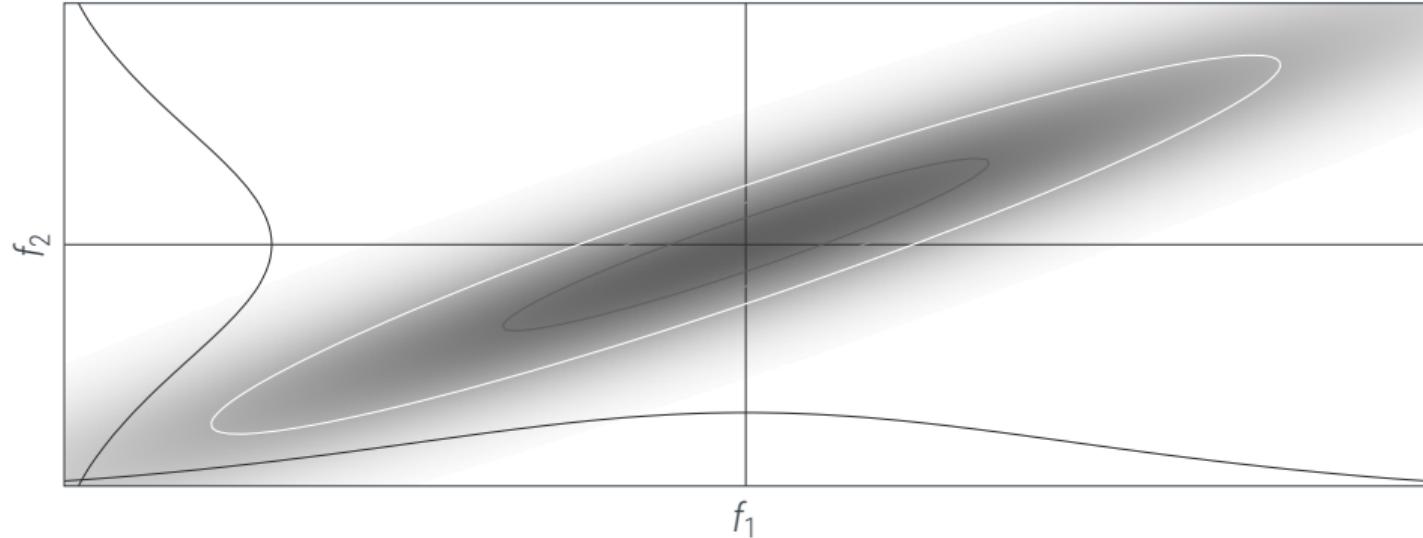
$$p(f_x | Y) = \frac{p(Y | f_x)p(f_x)}{p(Y)} = \frac{\sigma(Y \odot f_x)\mathcal{N}(f_x; m, k)}{\int \sigma(Y \odot f_x)\mathcal{N}(f_x; m, k) df_x}$$

$$\log p(f_x | Y) = -\frac{1}{2}f_x^\top k_{XX}^{-1}f_x + \sum_{i=1}^n \log \sigma(y_i f_{x_i}) + \text{const.}$$



# Logistic Regression is non-analytic

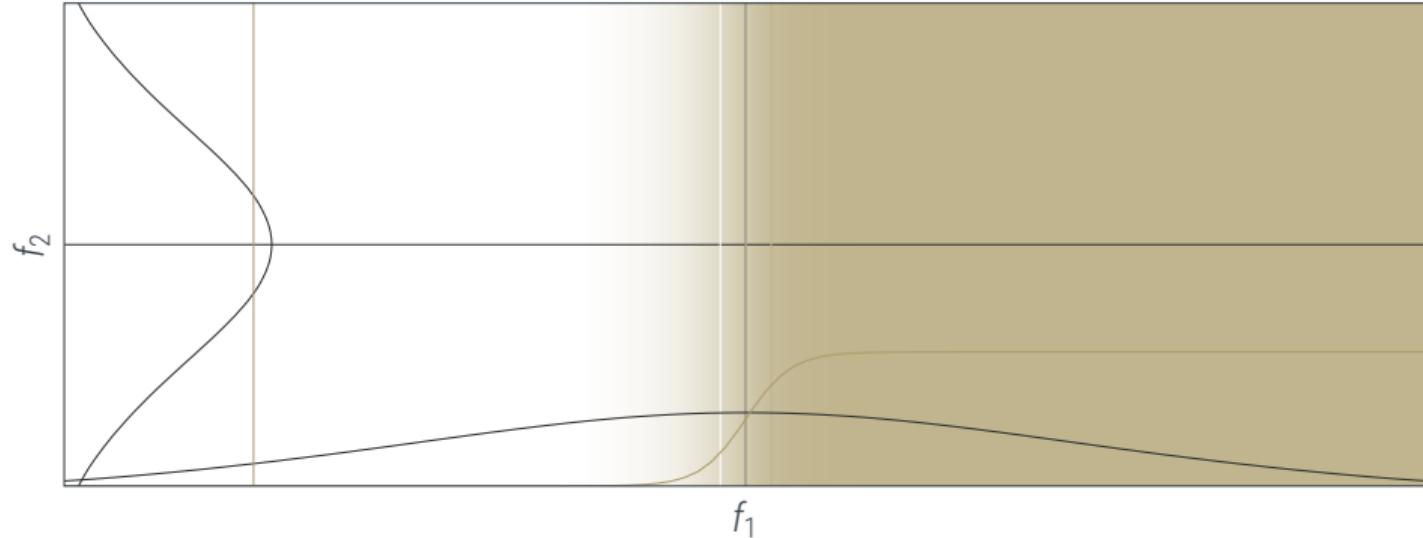
We'll have to break out the toolbox





# Logistic Regression is non-analytic

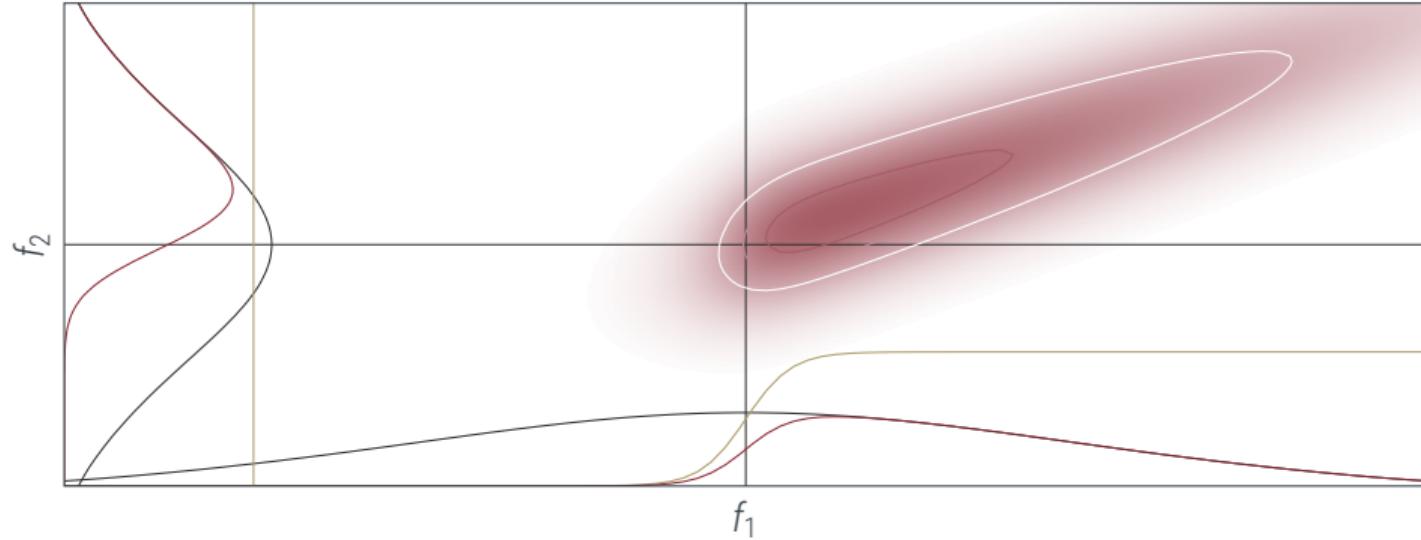
We'll have to break out the toolbox





# Logistic Regression is non-analytic

We'll have to break out the toolbox



We do not always care about all details of the posterior, just "key aspects".

- Remember that the Gaussian choice was also one of convenience.
- **Moments** of  $p(f, y) = p(y | f)p(f)$  we may be interested in

$$\mathbb{E}_p(1) = \int p(y, f) df = \int 1 \cdot p(y, f) df = Z \quad \text{the evidence}$$

$$\mathbb{E}_{p(f|y)}(f) = \int f \cdot p(f | y) df = \frac{1}{Z} \int f \cdot p(f, y) df = \bar{f} \quad \text{the mean}$$

$$\mathbb{E}_{p(f|y)}(f^2) - \bar{f}^2 = \int f^2 \cdot p(f | y) df - \bar{f}^2 = \frac{1}{Z} \int f^2 \cdot p(f, y) df - \bar{f}^2 = \text{var}(f) \quad \text{the variance}$$

$Z$  for hyperparameter tuning

$\bar{f}$  as a point estimator

$\text{var}(f)$  as an error estimator

Unfortunately, all these are usually intractable. But we can aim to approximate them.



## The Toolbox

Five principal methods for dealing with computational complexity in probabilistic inference

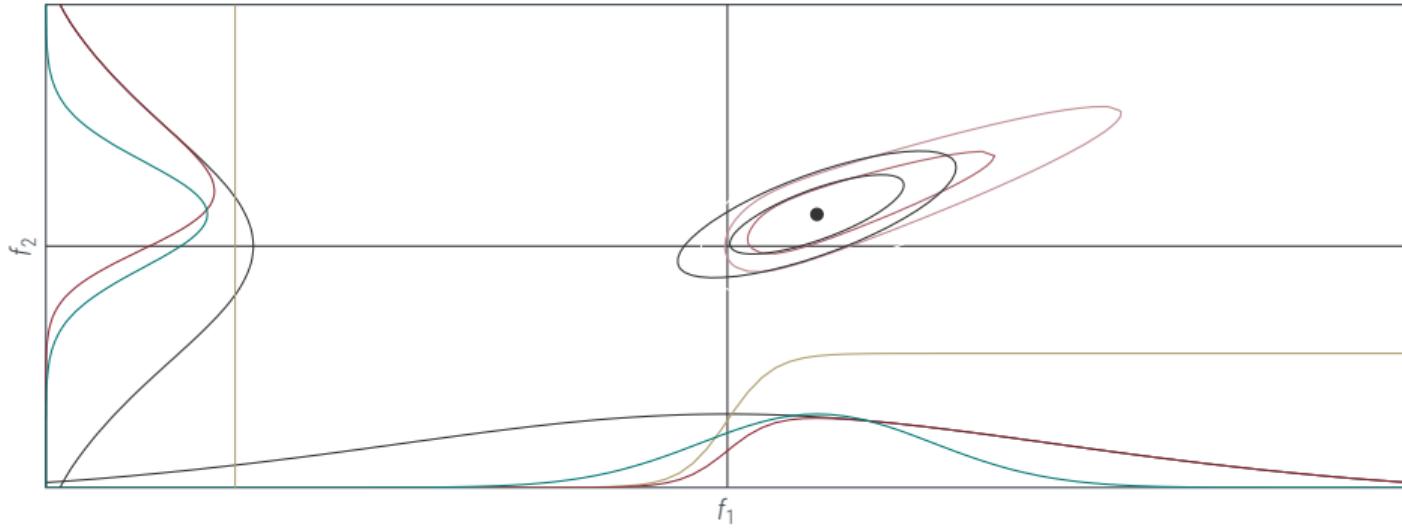
1. Maximum Likelihood (ML) / Maximum A-Posteriori (MAP) estimation:  
To estimate  $\theta$  in  $p(D | \theta)$  or  $p(\theta | D)$ , set  $\hat{\theta} = \arg \max_{\theta} p$ .
2. Laplace Approximation:  $p(\theta | D) \approx \mathcal{N}\left(\theta; \hat{\theta}, -(\nabla \nabla^T \log p(\theta | D))^{-1}\right)$
3. ????
4. ????
5. Numerical Quadrature:

To marginalize  $\theta$ , compute  $\int p(f | \theta) d\theta \approx \sum_i w_i \cdot p(f | \theta_i)$

Disclaimer: The listed items are neither mutually exclusive nor collectively exhaustive. Some of the methods are intricately interrelated.

# The Laplace Approximation

A local Gaussian approximation



# An idea as old as Probabilistic Inference

Théorie analytique des probabilités, 1814



Pierre-Simon, marquis de Laplace  
(1749-1827)

364

## THÉORIE ANALYTIQUE

L'intégrale du numérateur étant prise depuis  $x=0$  jusqu'à  $x=y$ , et celle du dénominateur étant prise depuis  $x=0$  jusqu'à  $x=a$

La valeur de  $x$  la plus probable, est celle qui rend  $y$  un *maximum*. Nous la désignerons par  $a$ . Si aux limites de  $x$ ,  $y$  est nul, alors chaque valeur de  $y$  a une valeur égale correspondante de l'autre côté du *maximum*.

Quand les valeurs de  $x$ , considérées indépendamment du résultat observé, ne sont pas également possibles; en nommant  $x$  la fonction de  $x$  qui exprime leur probabilité; il est facile de voir, par ce qui a été dit dans le premier chapitre de ce Livre, qu'en changeant dans la formule (1),  $y$  dans  $yz$ , on aura la probabilité que la valeur de  $x$  est comprise dans les limites  $x=\theta$  et  $x=\theta'$ . Cela revient à supposer toutes les valeurs de  $x$  également possibles *a priori*, et à considérer le résultat observé, comme étant formé de deux résultats indépendants, dont les probabilités sont  $y$  et  $z$ . On peut donc ramener ainsi tous les cas à celui où l'on suppose *a priori*, avant l'événement, une égale possibilité aux différentes valeurs de  $x$ , et par cette raison, nous adopterons cette hypothèse dans ce qui va suivre.

Nous avons donné dans les n° 22 et suivants du premier Livre, les formules nécessaires pour déterminer par des approximations convergentes, les intégrales du numérateur et du dénominateur de la formule (1), lorsque les événements simples dont se compose l'événement observé, sont répétés un très-grand nombre de fois; car alors  $y$  a pour facteurs, des fonctions de  $x$  élevées à de grandes puissances. Nous allons, au moyen de ces formules, déterminer la loi de probabilité des valeurs de  $x$ , à mesure qu'elles s'éloignent de la valeur  $a$ , la plus probable, ou qui rend  $y$  un *maximum*. Pour cela, reprenons la formule (c) du n° 27 du premier Livre,

$$\begin{aligned} \int y dx &= Y \cdot \left\{ U + \frac{1}{2} \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^2} + \frac{1 \cdot 3}{2!} \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot 4 \cdot dx^4} + \text{etc.} \right\} \cdot \int dt \cdot c^{-t} \cdot \\ &\quad + \frac{Y}{2} \cdot c^{-T} \cdot \left\{ \frac{d^2 U^2}{dx^2} - T \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^4} + (T+1) \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot dx^4} + \text{etc.} \right\}; \quad (2) \\ &\quad - \frac{Y}{2} \cdot c^{-T'} \cdot \left\{ \frac{d^2 U^2}{dx^2} + T' \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^4} + (T'+1) \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot dx^4} + \text{etc.} \right\}; \quad (3) \end{aligned}$$

## DES PROBABILITÉS.

565

$\psi$  est égal à  $\frac{x-a}{\sqrt{\log Y - \log y}}$ , et  $U$ ,  $\frac{d^2 U^2}{dx^2}$ ,  $\frac{d^3 U^3}{dx^4}$ , etc. sont ce qui deviennent  $v$ ,  $\frac{d^2 v^2}{dx^2}$ ,  $\frac{d^3 v^3}{dx^4}$ , etc., lorsqu'on  $y$  change après les différentiations,  $x$  en  $a$ ,  $a$  étant la valeur de  $x$  qui rend  $y$  un *maximum*:  $T$  est égal à ce que devient la fonction  $\sqrt{\log Y - \log y}$ , lorsqu'on change  $x$  en  $a-\theta$  dans  $y$ , et  $T'$  est ce que devient la même fonction, lorsqu'on  $y$  change  $x$  dans  $a+\theta$ . L'expression précédente de  $\int y dx$  donne la valeur de cette intégrale, dans les limites  $x=a-\theta$  et  $x=a+\theta$ ; l'intégrale  $\int dt \cdot c^{-t}$  étant prise depuis  $t=-T$  jusqu'à  $t=T'$ .

Le plus souvent, aux limites de l'intégrale  $\int y dx$ , étendue depuis  $x=0$  jusqu'à  $x=1$ ,  $y$  est nul; ou lorsque  $y$  n'est pas nul, il devient si petit à ces limites, qu'on peut le supposer nul. Alors, on peut faire à ces limites  $T$  et  $T'$  infinis, ce qui donne pour l'intégrale  $\int y dx$ , étendue depuis  $x=0$  jusqu'à  $x=1$ ,

$$\int y dx = Y \cdot \left\{ U + \frac{1}{2} \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^2} + \frac{1 \cdot 3}{2!} \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot 4 \cdot dx^4} + \text{etc.} \right\} \cdot \sqrt{\pi};$$

ainsi la probabilité que la valeur de  $x$  est comprise dans les limites  $x=a-\theta$  et  $x=a+\theta$ , est égale à

$$\frac{\int dt \cdot c^{-t} \cdot \left\{ \frac{1}{2} \cdot c^{-T} \cdot \left\{ \frac{d^2 U^2}{dx^2} - T \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^4} + (T+1) \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot dx^4} + \text{etc.} \right\} \right\}}{\sqrt{\pi}} + \frac{\int dt \cdot c^{-t} \cdot \left\{ \frac{1}{2} \cdot c^{-T'} \cdot \left\{ \frac{d^2 U^2}{dx^2} + T' \cdot \frac{d^2 U^2}{1 \cdot 2 \cdot dx^4} + (T'+1) \cdot \frac{d^3 U^3}{1 \cdot 2 \cdot 3 \cdot dx^4} + \text{etc.} \right\} \right\}}{\sqrt{\pi}}; \quad (3)$$

On voit par le n° 23 du premier Livre, que dans le cas où  $y$  a pour facteurs, des fonctions de  $x$  élevées à de grandes puissances de l'ordre  $\frac{1}{n}$ ,  $n$  étant une fraction extrêmement petite, alors  $U$  est le plus souvent de l'ordre  $\sqrt{n}$ , ainsi que ses différences successives;  $U$ ,  $\frac{d^2 U^2}{dx^2}$ ,  $\frac{d^3 U^3}{dx^4}$ , etc. sont respectivement des ordres  $\sqrt{n}$ ,  $n$ ,  $n^{\frac{1}{2}}$ , etc.; d'où il suit que la convergence des séries de la formule (3), exige que  $T$  et  $T'$  ne soient pas d'un ordre supérieur à  $\frac{1}{\sqrt{n}}$ .

# The Laplace Approximation

Once again, formally

- Consider a probability distribution  $p(\theta)$  (may be a posterior  $p(\theta | D)$  or something else)
- find a (local) **maximum** of  $p(\theta)$  or (equivalently)  $\log p(\theta)$

$$\hat{\theta} = \arg \max \log p(\theta) \quad \Rightarrow \quad \nabla \log p(\hat{\theta}) = 0$$

- perform **second order Taylor expansion** around  $\theta = \hat{\theta} + \delta$  in log space

$$\log p(\delta) = \log p(\hat{\theta}) + \frac{1}{2} \delta^\top \underbrace{\left( \nabla \nabla^\top \log p(\hat{\theta}) \right)}_{=: \Psi} \delta + \mathcal{O}(\delta^3)$$

- define **the Laplace approximation**  $q$  to  $p$

$$q(\theta) = \mathcal{N}(\theta; \hat{\theta}, -\Psi^{-1})$$

- Note that, if  $p(\theta) = \mathcal{N}(\theta; m, \Sigma)$ , then  $p(\theta) = q(\theta)$

# The Laplace Approximation for GP Classification

conceptual step (implementation details coming up)



[based on Rasmussen & Williams, 2006, §3.4]

- Find maximum posterior probability for **latent  $f$**  at **training points**

$$\hat{\mathbf{f}} = \arg \max \log p(\mathbf{f}_X | \mathbf{y})$$

- Assign approximate Gaussian posterior at training points

$$q(\mathbf{f}_X) = \mathcal{N}(\mathbf{f}_X; \hat{\mathbf{f}}, -(\nabla \nabla^T \log p(\mathbf{f}_X | \mathbf{y})|_{\mathbf{f}_X=\hat{\mathbf{f}}})^{-1}) =: \mathcal{N}(\mathbf{f}_X; \hat{\mathbf{f}}, \hat{\Sigma})$$

- approximate posterior **predictions** at  $f_x$  for **latent function**

$$\begin{aligned} q(f_x | \mathbf{y}) &= \int p(f_x | \mathbf{f}_X) q(\mathbf{f}_X) d\mathbf{f}_X = \int \mathcal{N}(f_x; m_x + k_{xx} K_{XX}^{-1} (\mathbf{f}_X - m_X), k_{xx} - k_{xx} K_{XX}^{-1} k_{xx}) q(\mathbf{f}_X) d\mathbf{f}_X \\ &= \mathcal{N}(f_x; m_x + k_{xx} K_{XX}^{-1} (\hat{\mathbf{f}} - m_X), k_{xx} - k_{xx} K_{XX}^{-1} k_{xx} + k_{xx} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{xx}) \end{aligned}$$

Compare with **exact predictions**

$$\mathbb{E}_{p(f_x, f_X | \mathbf{y})}(f_x) = \int (\mathbb{E}_{p(f_x | f_X)}(f_x)) p(f_X | \mathbf{y}) df_X = m_x + k_{xx} K_{XX}^{-1} (\mathbb{E}_{p(f_X | \mathbf{y})}(f_X) - m_X) =: \bar{f}_x$$

---

Recall:  $p(x) = \mathcal{N}(x; m, V)$ ,  $p(z | x) = \mathcal{N}(z; Ax, B) \Rightarrow p(z) = \int p(z | x)p(x) dx = \mathcal{N}(z; Am, AAV^T + B)$ .

# The Laplace Approximation for GP Classification

conceptual step (implementation details coming up)



[based on Rasmussen & Williams, 2006, §3.4]

- Find maximum posterior probability for **latent  $f$**  at **training points**

$$\hat{\mathbf{f}} = \arg \max \log p(\mathbf{f}_x | y)$$

- Assign approximate Gaussian posterior at training points

$$q(\mathbf{f}_x) = \mathcal{N}(\mathbf{f}_x; \hat{\mathbf{f}}, -(\nabla \nabla^T \log p(\mathbf{f}_x | y)|_{\mathbf{f}_x=\hat{\mathbf{f}}})^{-1}) =: \mathcal{N}(\mathbf{f}_x; \hat{\mathbf{f}}, \hat{\Sigma})$$

- approximate posterior **predictions** at  $f_x$  for **latent function**

$$\begin{aligned} q(f_x | y) &= \int p(f_x | \mathbf{f}_x) q(\mathbf{f}_x) d\mathbf{f}_x = \int \mathcal{N}(f_x; m_x + k_{xx} K_{xx}^{-1} (\mathbf{f}_x - m_x), k_{xx} - k_{xx} K_{xx}^{-1} k_{xx}) q(\mathbf{f}_x) d\mathbf{f}_x \\ &= \mathcal{N}(f_x; m_x + k_{xx} K_{xx}^{-1} (\hat{\mathbf{f}} - m_x), k_{xx} - k_{xx} K_{xx}^{-1} k_{xx} + k_{xx} K_{xx}^{-1} \hat{\Sigma} K_{xx}^{-1} k_{xx}) \end{aligned}$$

Compare with **exact predictions**

$$\text{var}_{p(f_x, f_x | y)}(f_x) = \int (f_x - \bar{f}_x)^2 dp(f_x | \mathbf{f}_x) dp(\mathbf{f}_x) = k_{xx} - k_{xx} K_{xx}^{-1} k_{xx} + k_{xx} K_{xx}^{-1} \text{var}_{p(f_x | y)}(f_x) K_{xx}^{-1} k_{xx}$$

---

Recall:  $p(x) = \mathcal{N}(x; m, V)$ ,  $p(z | x) = \mathcal{N}(z; Ax, B) \Rightarrow p(z) = \int p(z | x)p(x) dx = \mathcal{N}(z; Am, AAV^T + B)$ .

# The Laplace Approximation for GP Classification

conceptual step (implementation details coming up)



- Find maximum posterior probability for **latent  $f$**  at **training points**

$$\hat{\mathbf{f}} = \arg \max \log p(\mathbf{f}_x \mid \mathbf{y})$$

- Assign approximate Gaussian posterior at training points

$$q(\mathbf{f}_x) = \mathcal{N}(\mathbf{f}_x; \hat{\mathbf{f}}, -(\nabla \nabla^T \log p(\mathbf{f}_x \mid \mathbf{y})|_{\mathbf{f}_x=\hat{\mathbf{f}}})^{-1}) =: \mathcal{N}(\mathbf{f}_x; \hat{\mathbf{f}}, \hat{\Sigma})$$

- approximate posterior **predictions** at  $f_x$  for **latent function**

$$\begin{aligned} q(f_x \mid \mathbf{y}) &= \int p(f_x \mid \mathbf{f}_x) q(\mathbf{f}_x) d\mathbf{f}_x = \int \mathcal{N}(f_x; m_x + k_{xx} K_{xx}^{-1} (\mathbf{f}_x - m_x), k_{xx} - k_{xx} K_{xx}^{-1} k_{xx}) q(\mathbf{f}_x) d\mathbf{f}_x \\ &= \mathcal{N}(f_x; m_x + k_{xx} K_{xx}^{-1} (\hat{\mathbf{f}} - m_x), k_{xx} - k_{xx} K_{xx}^{-1} k_{xx} + k_{xx} K_{xx}^{-1} \hat{\Sigma} K_{xx}^{-1} k_{xx}) \end{aligned}$$

- compute predictions for **label probabilities**:

$$\mathbb{E}_{p(f \mid \mathbf{y})} [\pi_x] \approx \mathbb{E}_q [\pi_x] = \int \sigma(f_x) q(f_x \mid \mathbf{y}) df_x \quad \text{or (not the same!) } \hat{\pi}_x = \sigma(\mathbb{E}_q(f_x))$$



- the Laplace approximation is only very roughly motivated (see above)
- it can be **arbitrarily wrong**, since it is a **local** approximation
- but it is often among the most computationally efficient things to try
- for logistic regression, it tends to work relatively well, because
  - the log posterior is concave (see below)
  - the algebraic structure of the link function yields "almost" a Gaussian posterior (cf. picture above)

# Implementing the Laplace Approximation

Derivation



[based on Rasmussen & Williams, 2006, §3.4]

$$p(f) = \mathcal{GP}(f, m, k) \quad p(\mathbf{y} | f_X) = \sigma(\mathbf{y} \odot f_X)$$

$$\log p(f_X | \mathbf{y}) = \log p(\mathbf{y} | f_X) + \log p(f_X) - \log p(\mathbf{y})$$

$$= \sum_{i=1}^n \log \sigma(y_i f_{x_i}) - \frac{1}{2} (\mathbf{f}_X - \mathbf{m}_X)^T K_{XX}^{-1} (\mathbf{f}_X - \mathbf{m}_X) + \text{const.}$$

$$\nabla \log p(f_X | \mathbf{y}) = \sum_{i=1}^n \nabla \log \sigma(y_i f_{x_i}) - K_{XX}^{-1} (\mathbf{f}_X - \mathbf{m}_X) \quad \text{with} \quad \frac{\partial \log \sigma(y_i f_{x_i})}{\partial f_{x_j}} = \delta_{ij} \left( \frac{y_i + 1}{2} - \sigma(f_{x_i}) \right)$$

$$\begin{aligned} \nabla \nabla^T \log p(f_X | \mathbf{y}) &= \sum_{i=1}^n \nabla \nabla^T \log \sigma(y_i f_{x_i}) - K_{XX}^{-1} \quad \text{with} \quad \frac{\partial^2 \log \sigma(y_i f_{x_i})}{\partial f_{x_a} \partial f_{x_b}} = -\delta_{ia} \delta_{ib} \underbrace{\sigma(f_{x_i})(1 - \sigma(f_{x_i}))}_{=: w_i \text{ with } 0 < w_i < 1} \\ &=: -\text{diag } \mathbf{w} - K^{-1} = -(W + K^{-1}) \quad \leftarrow \text{convex minimization / concave maximization!} \end{aligned}$$

# Implementing the Laplace Approximation I

Newton Optimization



```
1 procedure OPTIMIZE( $L(\cdot)$ ,  $f_0$ )
2    $f \leftarrow f_0$                                      // initialize
3   while not converged do
4      $g \leftarrow \nabla L(f)$                          // compute gradient
5      $H \leftarrow (\nabla \nabla^\top L(f))^{-1}$         // compute inverse Hessian
6      $\Delta \leftarrow Hg$                             // Newton step
7      $f \leftarrow f - \Delta$                          // perform step
8     converged  $\leftarrow \|\Delta\| < \epsilon$           // check for convergence
9   end while
10  return  $f$ 
11 end procedure
```

# Implementing the Laplace Approximation II

Concrete Algorithm, preliminary form



```
1 procedure GP-LOGISTIC-TRAIN( $K_{XX}$ ,  $m_X$ ,  $y$ )
2    $f \leftarrow m_X$  // initialize
3   while not converged do
4      $r \leftarrow \frac{y+1}{2} - \sigma(f)$  // =  $\nabla \log p(y | f_X)$ , gradient of log likelihood
5      $W \leftarrow \text{diag}(\sigma(f) \odot (1 - \sigma(f)))$  // =  $-\nabla \nabla \log p(y | f_X)$ , Hessian of log likelihood
6      $g \leftarrow r - K_{XX}^{-1}(f - m_X)$  // compute gradient
7      $H \leftarrow -(W + K^{-1})^{-1}$  // compute inverse Hessian
8      $\Delta \leftarrow Hg$  // Newton step
9      $f \leftarrow f - \Delta$  // perform step
10    converged  $\leftarrow \|\Delta\| < \epsilon$  // check for convergence
11  end while
12  return  $f$ 
13 end procedure
```

This can be numerically unstable as it (repeatedly) requires  $(W + K^{-1})^{-1}$ . For a numerically stable alternative, see end of this pdf.



# Computing Predictions

from the Laplace approximation

[Rasmussen & Williams, 2006, §3.4]

$$\log p(f_X | y) = \sum_{i=1}^n \log \sigma(y_i f_{x_i}) - \frac{1}{2}(f_X - m_X)^T K_{XX}^{-1}(f_X - m_X) + \text{const.}$$

$$\nabla \log p(f_X | y) = \underbrace{\sum_{i=1}^n \nabla \log \sigma(y_i f_{x_i})}_{=: r} - K_{XX}^{-1}(f_X - m_X) \quad \text{with} \quad \frac{\partial \log \sigma(y_i f_{x_i})}{\partial f_{x_j}} = \delta_{ij} \left( \frac{y_i + 1}{2} - \sigma(f_{x_i}) \right)$$

```

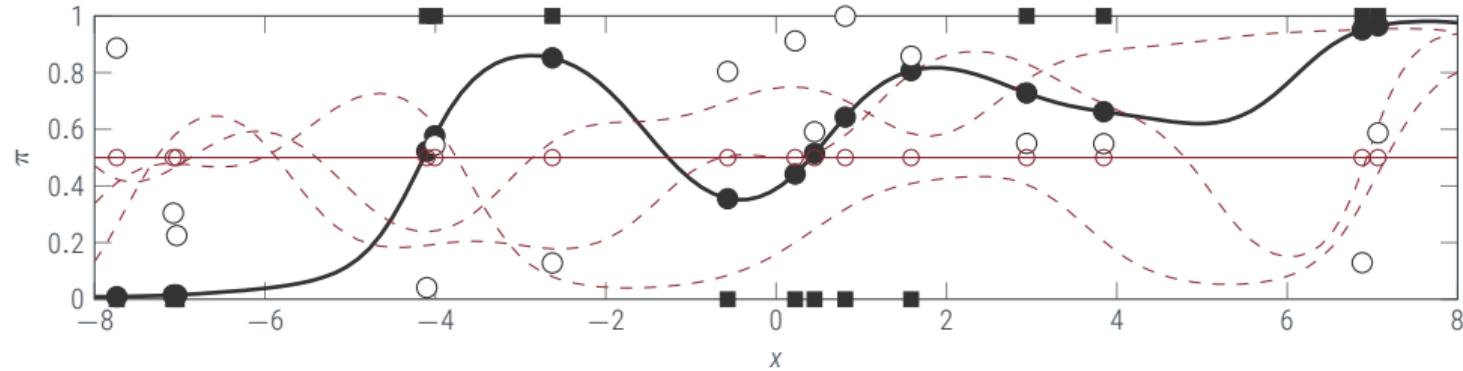
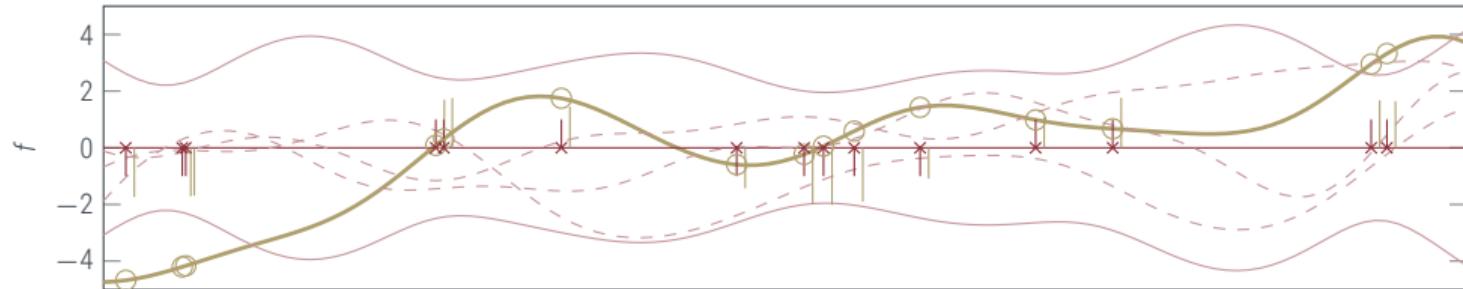
1 procedure GP-LOGISTIC-PREDICT( $\hat{f}, W, R, r, k, x$ ) //  $\hat{f}, W, R, r$  handed over from training
2   for  $i = 1, \dots, \text{LENGTH}(x)$  do
3      $\bar{f}_i \leftarrow k_{x_i X} r$  // mean prediction (note at minimum,  $0 = \nabla p(f_X | y) = r - K_{XX}^{-1}(f_X - m_X)$ ).
4      $s \leftarrow R^{-1}(W^{1/2}k_{X x_i})$  // pre-computation allows this step in  $\mathcal{O}(n^2)$ 
5      $v \leftarrow k_{x_i X_i} - s^T s$  //  $v = \text{cov}(f_x)$ 
6      $\bar{\pi}_i \leftarrow \int \sigma(\bar{f}_i) \mathcal{N}(f_i, \bar{f}_i, v) df_i$  // predictive probability for class 1 is  $p(y | \bar{f}) = \int p(y_x | f_x)p(f_x | \bar{f}) df_x$ 
7   end for // entire loop is  $\mathcal{O}(n^2m)$  for  $m$  test cases.
8   return  $\bar{\pi}_x$ 
9 end procedure

```



# Pictorial View

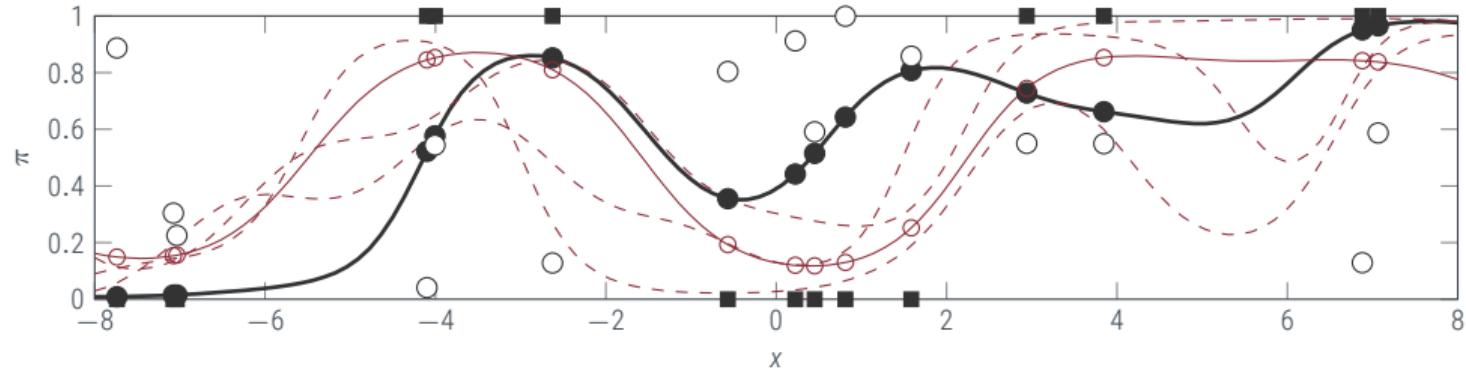
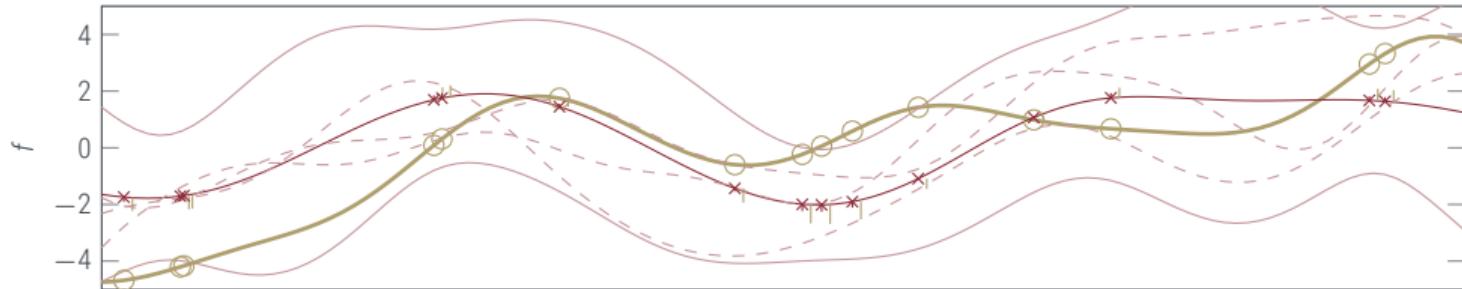
Gaussian process logistic regression





# Pictorial View

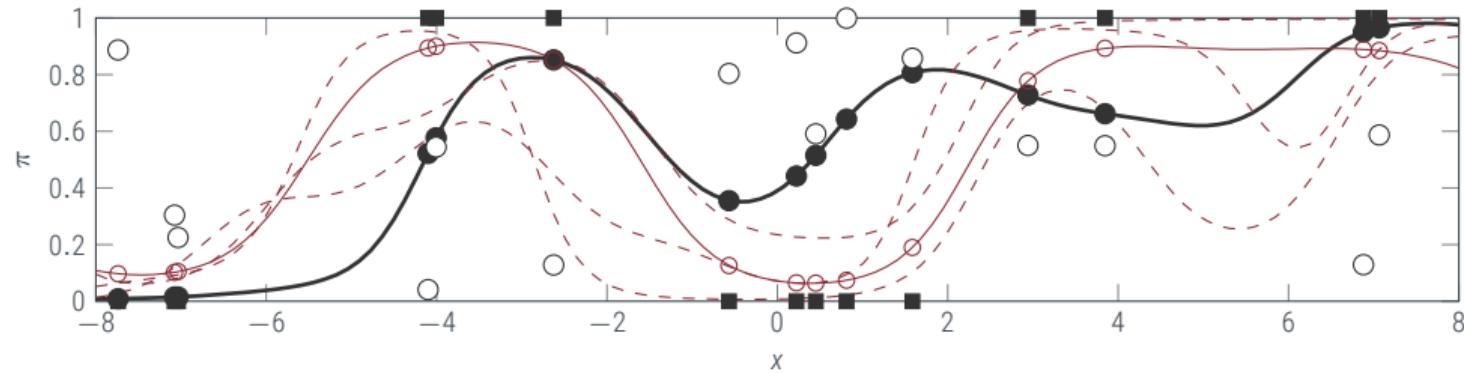
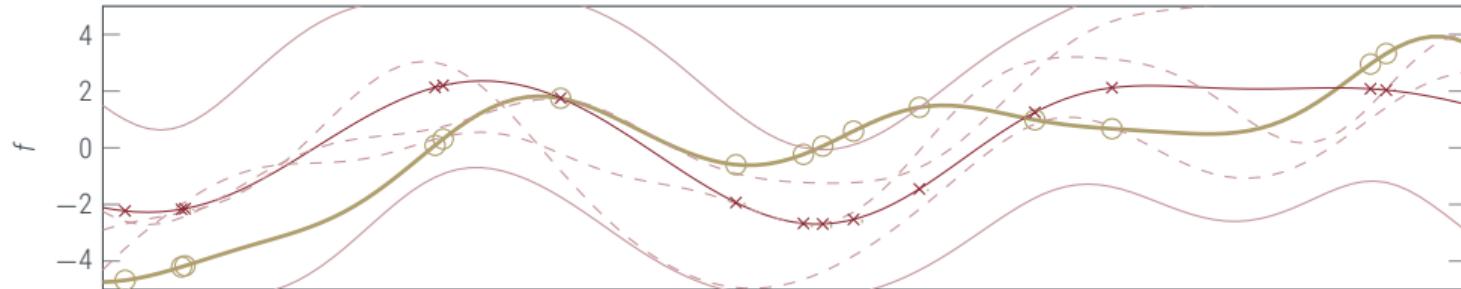
Gaussian process logistic regression





# Pictorial View

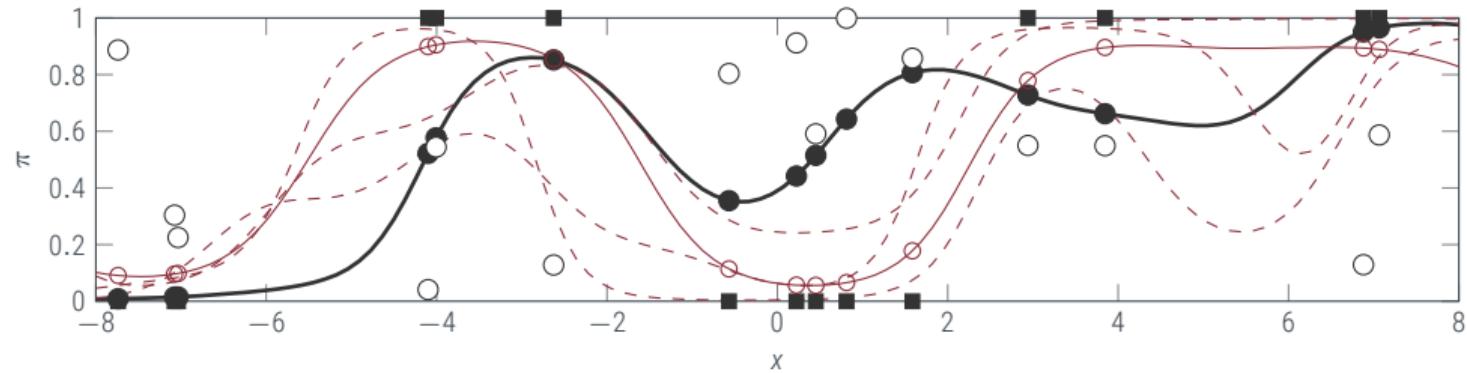
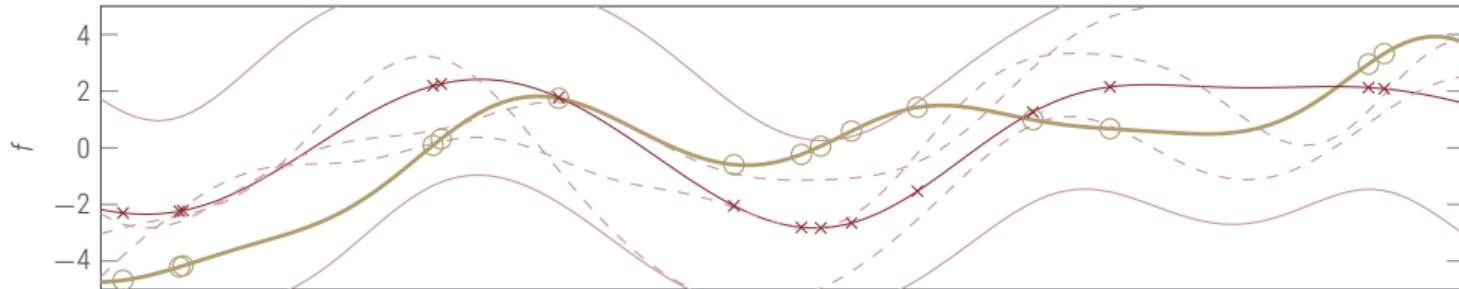
Gaussian process logistic regression





# Pictorial View

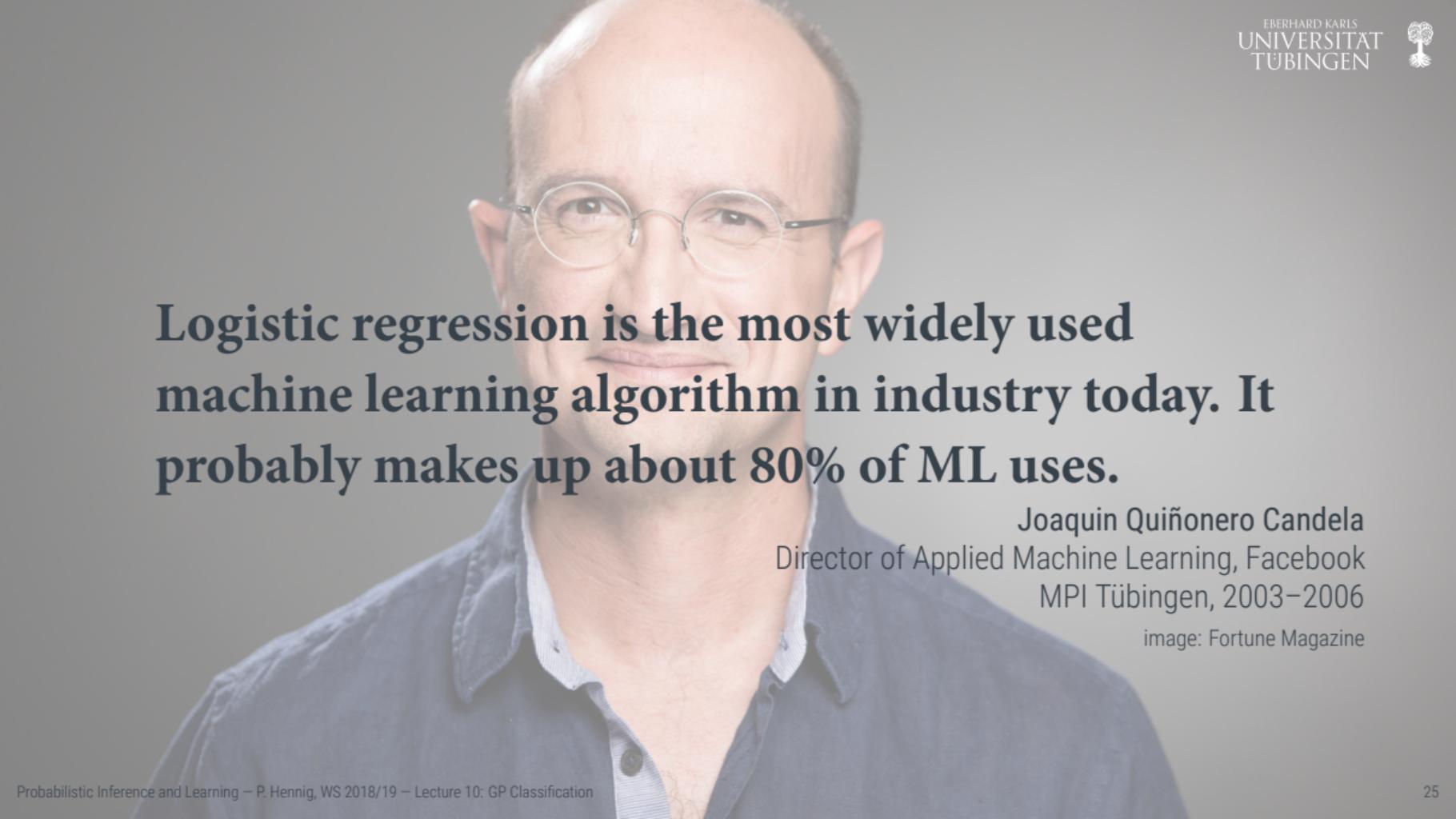
Gaussian process logistic regression





## Gaussian Process Classification – (Probabilistic) Logistic Regression:

- Supervised classification phrased in a **discriminative** model with probabilistic interpretation
- model binary outputs as a **transformation** of a **latent function** with a Gaussian process prior
- due to **non-Gaussian likelihood**, the posterior is non-Gaussian; exact inference **intractable**
- **Laplace approximation**: Find MAP estimator, second order expansion for Gaussian approximation
- tune code for numerical stability, efficient computations
- Laplace approximation provides Gaussian posterior on training points, hence evidence, predictions



**Logistic regression is the most widely used  
machine learning algorithm in industry today. It  
probably makes up about 80% of ML uses.**

Joaquin Quiñonero Candela  
Director of Applied Machine Learning, Facebook  
MPI Tübingen, 2003–2006

image: Fortune Magazine

# Evidence / Marginal Likelihood

For hyperparameter adaptation



$$p(\mathbf{y} | X) = \int p(\mathbf{y}, \mathbf{f} | X) d\mathbf{f} = \int \exp(\log p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | X)) d\mathbf{f}$$

Laplace:  $\log p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | X) \approx \log(p(\mathbf{y} | \hat{\mathbf{f}}) p(\hat{\mathbf{f}} | X)) - \frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^T(K^{-1} + W)(\mathbf{f} - \hat{\mathbf{f}}) = q(\mathbf{y}, \mathbf{f} | X)$

thus  $p(\mathbf{y} | X) \approx q(\mathbf{y} | X) = \exp\left(\log\left(p(\mathbf{y} | \hat{\mathbf{f}}) p(\hat{\mathbf{f}} | X)\right)\right) \int \exp\left(-\frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^T(K^{-1} + W)(\mathbf{f} - \hat{\mathbf{f}})\right) d\mathbf{f}$

$$= \exp\left(\log\left(p(\mathbf{y} | \hat{\mathbf{f}})\right)\right) \mathcal{N}(\hat{\mathbf{f}}; m_X, k_{XX})(2\pi)^{n/2}|(K^{-1} + W)^{-1}|^{1/2}$$

$$\log q(\mathbf{y} | X) = \log p(\mathbf{y} | \mathbf{f}) - \frac{1}{2}(\hat{\mathbf{f}} - m_X)^T K_{XX}^{-1}(\hat{\mathbf{f}} - m_X) - \frac{1}{2} \log(|K| \cdot |K^{-1} + W|)$$

$$= \sum_{i=1}^n \sigma(y_i f_{x_i}) - \frac{1}{2}(\hat{\mathbf{f}} - m_X)^T K_{XX}^{-1}(\hat{\mathbf{f}} - m_X) - \frac{1}{2} \log |B|$$

# Evidence / Marginal Likelihood

For hyperparameter adaptation



```
1 procedure GP-LOGISTIC-TRAIN( $K_{XX}$ ,  $m_X$ ,  $y$ )
2    $f \leftarrow m_X$  // initialize
3   while not converged do
4      $r \leftarrow \frac{y+1}{2} - \sigma(f)$  //  $= \nabla \log p(y | f_X)$ , gradient of log likelihood
5      $W \leftarrow \text{diag}(\sigma(f) \odot (1 - \sigma(f)))$  //  $= -\nabla \nabla \log p(y | f_X)$ , Hessian of log likelihood
6      $R \leftarrow \text{CHOLESKY}(I + W^{1/2}K_{XX}W^{1/2})$  //  $B$  from previous slide
7      $b \leftarrow W(f - m_X) + r$ 
8      $a \leftarrow b - W^{1/2}R^{-\top}R^{-1}(W^{1/2}Kb)$  //  $a = K^{-1}(K^{-1} + W)^{-1}b$ 
9      $f \leftarrow Ka + m_X$  //  $f \leftarrow (K^{-1} + W)^{-1}((K^{-1} + W)(f - m_X) - K^{-1}(f - m_X) + r + m_X) = f + (W + K^{-1})^{-1}g$ 
10    end while // The objective is  $-\frac{1}{2}a^\top(f - m_X) + \sum_i \log \sigma(y_i f_i)$ 
11     $\log q(y | X) \leftarrow -\frac{1}{2}a^\top(f - m_X) + \sum_i \log \sigma(y_i f_i) - \sum_i \log R_{ii}$  // (Cholesky  $B = R^\top R \Rightarrow |B| = \prod_i R_{ii}^2$ )
12    return  $f, W, R, r$ 
13 end procedure
```

# Parametric Form

For completeness – analogous to above

$$p(v) = \mathcal{N}(v, \mu, \Sigma) \quad p(y | v) = \sigma(y \odot \phi_X^\top v) \quad \text{with } v \in \mathbb{R}^F, \phi_X \in \mathbb{R}^{n \times F}$$

$$\log p(v | y) = \log p(y | v) + \log p(v) - \log p(y)$$

$$= \sum_{i=1}^n \log \sigma(y_i \phi_{x_i}^\top v) - \frac{1}{2}(v - \mu)^\top \Sigma^{-1}(v - \mu) + \text{const.}$$

$$\nabla \log p(v | y) = \sum_{i=1}^n \nabla \log \sigma(y_i \phi_{x_i}^\top v) - \Sigma^{-1}(v - \mu) \quad \text{with} \quad \frac{\partial \log \sigma(y_i \phi_{x_i}^\top v)}{\partial v_j} = [\phi_{x_i}]_j \left( \frac{y_i + 1}{2} - \sigma(\phi_{x_i}^\top v) \right)$$

$$\nabla \nabla^\top \log p(v | y) = \sum_{i=1}^n \nabla \nabla^\top \log \sigma(y_i \phi_{x_i}^\top v) - \Sigma^{-1} \quad \text{with} \quad \frac{\partial^2 \log \sigma(y_i \phi_{x_i}^\top v)}{\partial v_a \partial v_b} = -[\phi_{x_i}]_a [\phi_{x_i}]_b \underbrace{\sigma(\phi_{x_i}^\top v)(1 - \sigma(\phi_{x_i}^\top v))}_{=: w_i}$$

$$=: -(W + \Sigma^{-1}) = -(\phi_X \text{diag}(w) \phi_X^\top + \Sigma^{-1}) \in \mathbb{R}^{F \times F}$$

still convex minimization / concave maximization! All computations  $\mathcal{O}(F^2 n)$

Exercise: Re-write Algorithms GP-Logistic-Train & GP-Logistic-Predict in weight-space.



# Multi-Class Classification

Extending the Laplace Approximation

- What if we have  $C$  class labels  $[c_1, \dots, c_C]$ ?



# Multi-Class Classification

Extending the Laplace Approximation

- What if we have  $C$  class labels  $[c_1, \dots, c_C]$ ?
- Generative Model:



# Multi-Class Classification

Extending the Laplace Approximation

- What if we have  $C$  class labels  $[c_1, \dots, c_C]$ ?
- Generative Model:
  - use  $C$  outputs of the latent function. So at the  $n$  locations, the latent variables are

$$\mathbf{f}_x = [f_1^{(1)}, \dots, f_n^{(1)}, f_1^{(2)}, \dots, f_n^{(2)}, \dots, f_1^{(C)}, \dots, f_n^{(C)}]$$

# Multi-Class Classification

Extending the Laplace Approximation

- What if we have  $C$  class labels  $[c_1, \dots, c_C]$ ?
- Generative Model:
  - use  $C$  outputs of the latent function. So at the  $n$  locations, the latent variables are

$$\mathbf{f}_x = [f_1^{(1)}, \dots, f_n^{(1)}, f_1^{(2)}, \dots, f_n^{(2)}, \dots, f_1^{(C)}, \dots, f_n^{(C)}]$$

- At location  $x_i$ , generate probabilities for each class by taking the **softmax**

$$p(y_i^{(c)} \mid \mathbf{f}_i) = \pi_i^{(c)} = \frac{\exp(f_i^{(c)})}{\sum_{\tilde{c}=1}^C \exp(f_i^{(\tilde{c})})}$$



# Multi-Class Classification

Extending the Laplace Approximation

- What if we have  $C$  class labels  $[c_1, \dots, c_C]$ ?
- Generative Model:
  - use  $C$  outputs of the latent function. So at the  $n$  locations, the latent variables are

$$\mathbf{f}_x = [f_1^{(1)}, \dots, f_n^{(1)}, f_1^{(2)}, \dots, f_n^{(2)}, \dots, f_1^{(C)}, \dots, f_n^{(C)}]$$

- At location  $x_i$ , generate probabilities for each class by taking the **softmax**

$$p(y_i^{(c)} \mid \mathbf{f}_i) = \pi_i^{(c)} = \frac{\exp(f_i^{(c)})}{\sum_{\tilde{c}=1}^C \exp(f_i^{(\tilde{c})})}$$

The remaining derivations are analogous to the binary case. They are this week's **exercise**.



### Definition (Generalized Linear Model)

(For our purposes,) a **generalized linear model (GLM)** is a probabilistic regression model for a function  $f$  with a Gaussian process prior  $p(f)$  and a **non-Gaussian** likelihood  $p(y | f_x)$ . Note the distinction to a **general** linear model (GP prior and likelihood, with non-linear kernel  $k$ )



## Summary:

- **Gaussian process classification** provides a probabilistic, discriminative model for classification
- The MAP solution is also known as **logistic regression**, and may be the most widely used machine learning algorithm in industry
- It can be found, for example by **Newton optimization**, or by stochastic gradient descent. The **Laplace approximation** provides a simple approximate posterior (it requires evaluating or estimating the *Hessian* of the log posterior)
- Other non-linearities / link-functions can be used, too, to model other output types (bounded, strictly positive, structured, etc.)

next time: Concrete example implementation

# Some practical numerical linear algebra

no, this will not be in the exam ...



The matrix  $(W + K^{-1})^{-1}$  can be numerically unstable. Instead, consider

$$B := I + W^{1/2}K_{XX}W^{1/2}.$$

- can be computed in  $\mathcal{O}(n^2)$  (because  $W$  is diagonal, spd)
- eigenvalues  $\lambda_i(B)$  are bounded  $1 < \lambda_i < 1 + n \max_{ij} \frac{K_{ij}}{4}$ , thus typically numerically stable
- Matrix inversion lemma yields

$$(K^{-1} + W)^{-1} = K - KW^{1/2}B^{-1}W^{1/2}K$$

$$(W^{-1} + K)^{-1} = W^{1/2}W^{-1/2}(W^{-1} + K)^{-1}W^{-1/2}W^{1/2} = W^{1/2}B^{-1}W^{1/2}$$

Lemma (The matrix inversion lemma / Sherman-Morrison-Woodbury Identity)

If all the inverses exist, then

$$(A^{-1} + UD^{-1}V^\top)^{-1} = A - AU(D + V^\top AU)^{-1}V^\top A.$$

# Some practical numerical linear algebra

no, this will not be in the exam ...



The matrix  $(W + K^{-1})^{-1}$  can be numerically unstable. Instead, consider

$$B := I + W^{1/2}K_{XX}W^{1/2}.$$

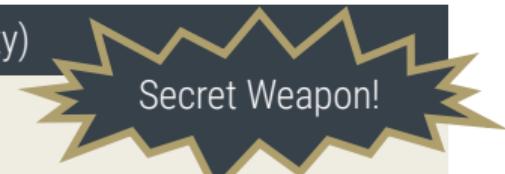
- can be computed in  $\mathcal{O}(n^2)$  (because  $W$  is diagonal, spd)
- eigenvalues  $\lambda_i(B)$  are bounded  $1 < \lambda_i < 1 + n \max_{ij} \frac{K_{ij}}{4}$ , thus typically numerically stable
- Matrix inversion lemma yields

$$(K^{-1} + W)^{-1} = K - KW^{1/2}B^{-1}W^{1/2}K$$

$$(W^{-1} + K)^{-1} = W^{1/2}W^{-1/2}(W^{-1} + K)^{-1}W^{-1/2}W^{1/2} = W^{1/2}B^{-1}W^{1/2}$$

Lemma (The matrix inversion lemma / Sherman-Morrison-Woodbury Identity)

If all the inverses exist, then



$$(A^{-1} + UD^{-1}V^\top)^{-1} = A - AU(D + V^\top AU)^{-1}V^\top A.$$

# Implementing the Laplace Approximation III

Efficient form



```
1 procedure GP-LOGISTIC-TRAIN( $K_{XX}$ ,  $m_X$ ,  $y$ )
2    $f \leftarrow m_X$  // initialize
3   while not converged do
4      $r \leftarrow \frac{y+1}{2} - \sigma(f)$  //  $= \nabla \log p(y | f_X)$ , gradient of log likelihood
5      $W \leftarrow \text{diag}(\sigma(f) \odot (1 - \sigma(f)))$  //  $= -\nabla \nabla \log p(y | f_X)$ , Hessian of log likelihood
6      $R \leftarrow \text{CHOLESKY}(I + W^{1/2}K_{XX}W^{1/2})$  //  $B$  from previous slide
7      $b \leftarrow W(f - m_X) + r$ 
8      $a \leftarrow b - W^{1/2}R^{-\top}R^{-1}(W^{1/2}Kb)$  //  $a = K^{-1}(K^{-1} + W)^{-1}b$ 
9      $f \leftarrow Ka + m_X$  //  $f \leftarrow (K^{-1} + W)^{-1}((K^{-1} + W)(f - m_X) - K^{-1}(f - m_X) + r + m_X) = f + (W + K^{-1})^{-1}g$ 
10   end while // The objective is  $-\frac{1}{2}a^\top f + \sum_i \log \sigma(y_i f_i)$ 
11   return  $f, W, R, r$ 
12 end procedure
```