# PROBABILISTIC MACHINE LEARNING LECTURE 15 GAUSSIAN PROCESS CLASSIFICATION

Philipp Hennig 22 June 2023

UNIVERSITÄT TÜBINGEN



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

## A Gaussian Process model for Classification



Logistic Regression

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

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

## A Gaussian Process model for Classification



Logistic Regression

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

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

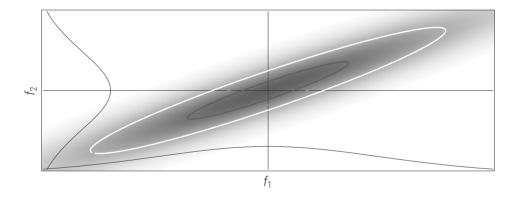
The problem: The posterior is not Gaussian!

$$p(f_X \mid Y) = \frac{p(Y \mid f_X)p(f_X)}{p(Y)} = \frac{\mathcal{N}(f_X; m, k) \prod_{i=1}^n \sigma(y_i f_{x_i})}{\int \mathcal{N}(f_X; m, k) \prod_{i=1}^n \sigma(y_i f_{x_i}) df_X}$$
$$\log p(f_X \mid Y) = -\frac{1}{2} f_X^\mathsf{T} k_{XX}^{-1} f_X + \sum_{i=1}^n \log \sigma(y_i f_{x_i}) + \text{const.}$$

# Logistic Regression is not analytic

eberhard karls UNIVERSITAT TUBINGEN

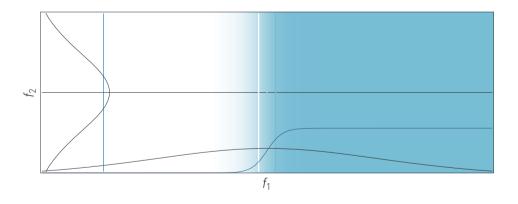
We'll have to break out the toolbox



# Logistic Regression is not analytic



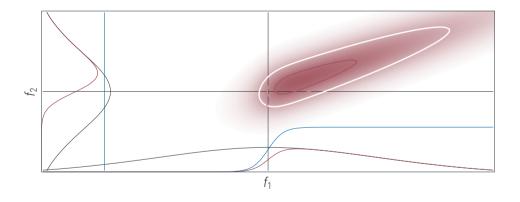
We'll have to break out the toolbox



# Logistic Regression is not analytic



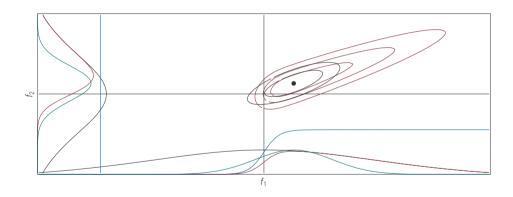
We'll have to break out the toolbox



# The Laplace Approximation

A local Gaussian approximation

erre Simon M. de Laplace, 1814



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

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

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

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

define the Laplace approximation q to p

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

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

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

Assign approximate Gaussian posterior at training points

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

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

$$q(f_X \mid y) = \int p(f_X \mid f_X)q(f_X) df_X = \int \mathcal{N}(f_X; m_X + k_{XX}K_{XX}^{-1}(f_X - m_X), k_{XX} - k_{XX}K_{XX}^{-1}k_{XX})q(f_X) df_X$$
  
=  $\mathcal{N}(f_X; m_X + k_{XX}K_{XX}^{-1}(\hat{\mathbf{f}} - m_X), k_{XX} - k_{XX}K_{XX}^{-1}k_{XX} + k_{XX}K_{XX}^{-1}\hat{\Sigma}K_{XX}^{-1}k_{XX})$ 

Compare with exact predictions

$$\mathbb{E}_{p(f_{x},f_{X}|y)}(f_{x}) = \int (\mathbb{E}_{p(f_{x}|f_{X})}(f_{x}))p(f_{X}|y) df_{X} = m_{X} + k_{XX}K_{XX}^{-1}(\mathbb{E}_{p(f_{X}|y)}(f_{X}) - m_{X}) =: \bar{f}_{X}$$

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



conceptual step (implementation details coming up

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

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

Assign approximate Gaussian posterior at training points

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

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

$$q(f_X \mid y) = \int p(f_X \mid f_X) q(f_X) df_X = \int \mathcal{N}(f_X; m_X + k_{XX} K_{XX}^{-1} (f_X - m_X), k_{XX} - k_{XX} K_{XX}^{-1} k_{XX}) q(f_X) df_X$$

$$= \mathcal{N}(f_X; m_X + k_{XX} K_{XX}^{-1} (\hat{f} - m_X), k_{XX} - k_{XX} K_{XX}^{-1} k_{XX} + k_{XX} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{XX})$$

Compare with exact predictions

$$\operatorname{var}_{p(f_{x},f_{X}\mid y)}(f_{x}) = \int (f_{x} - \bar{f}_{x})^{2} dp(f_{x}\mid f_{x}) dp(f_{x}) = k_{xx} - k_{xx}K_{xx}^{-1}k_{xx} + k_{xx}K_{xx}^{-1}\operatorname{var}_{p(f_{x}\mid y)}(f_{x})K_{xx}^{-1}k_{xx}$$

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



# The Laplace Approximation for GP Classification

conceptual step (implementation details coming up

ased on Rasmussen & Williams, 2006, §3.4]

► Find maximum posterior probability for **latent** *f* at **training points** 

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

► Assign approximate Gaussian posterior at training points

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

 $\triangleright$  approximate posterior **predictions** at  $f_x$  for **latent function** 

$$q(f_X \mid y) = \int p(f_X \mid f_X) q(f_X) df_X = \int \mathcal{N}(f_X; m_X + k_{XX} K_{XX}^{-1} (f_X - m_X), k_{XX} - k_{XX} K_{XX}^{-1} k_{XX}) q(f_X) df_X$$

$$= \mathcal{N}(f_X; m_X + k_{XX} K_{XX}^{-1} (\hat{f} - m_X), k_{XX} - k_{XX} K_{XX}^{-1} k_{XX} + k_{XX} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{XX})$$

compute predictions for label probabilities:

$$\mathbb{E}_{p(f|y)}[\pi_X] \approx \mathbb{E}_q[\pi_X] = \int \sigma(f_X) q(f_X \mid y) \, df_X \quad \text{or (not the same!)} \quad \hat{\pi}_X = \sigma(\mathbb{E}_q(f_X))$$



- the Laplace approximation is only very roughly motivated (see above)
- it can be **arbitrarily wrong**, since it is a **local** approximation
- but it is still better than a point estimate!
- and it is typically the most computationally efficient thing to try, because it uses only auto-diff and linear algebra
- 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),

Today, we will focus on the Hessian/uncertainty

ased on Rasmussen & Williams, 2006, §3.4]

$$p(f) = \mathcal{GP}(f, m, k) \qquad p(\mathbf{y} \mid f_{X}) = \prod_{i=1}^{n} \sigma(y_{i} f_{x_{i}}) \qquad \sigma(z) = \frac{1}{1 + e^{-x}}$$

$$\log p(f_{X} \mid \mathbf{y}) = \log p(\mathbf{y} \mid f_{X}) + \log p(f_{X}) - \log p(\mathbf{y}) \quad \text{with} \quad \log \sigma(y_{i} f_{x_{i}}) = -\log(1 + e^{-y_{i} f_{x_{i}}})$$

$$= \sum_{i=1}^{n} \log \sigma(y_{i} f_{x_{i}}) - \frac{1}{2} (f_{X} - m_{X})^{\mathsf{T}} K_{XX}^{-1} (f_{X} - m_{X}) + \text{const.}$$

$$\nabla \log p(f_{X} \mid \mathbf{y}) = \sum_{i=1}^{n} \nabla \log \sigma(y_{i} f_{x_{i}}) - K_{XX}^{-1} (f_{X} - m_{X}) \quad \text{with} \quad \frac{\partial \log \sigma(y_{i} f_{x_{i}})}{\partial f_{x_{i}}} = \delta_{ij} \left( \frac{y_{i} + 1}{2} - \sigma(f_{x_{i}}) \right)$$

$$\nabla \nabla^{\mathsf{T}} \log p(f_{X} \mid \mathbf{y}) = \sum_{i=1}^{n} \nabla \nabla^{\mathsf{T}} \log \sigma(y_{i} f_{x_{i}}) - K_{XX}^{-1} \quad \text{with} \quad \frac{\partial^{2} \log \sigma(y_{i} f_{x_{i}})}{\partial f_{x_{a}} \partial f_{x_{b}}} = -\delta_{ia} \delta_{ib} \underbrace{\sigma(f_{x_{i}}) (1 - \sigma(f_{x_{i}}))}_{=:w_{i} \text{ with } 0 < w_{i} < 1}$$

$$=: -\operatorname{diag} \mathbf{w} - K^{-1} = -(W + K^{-1}) \quad \leftarrow \text{convex minimization} / \text{concave maximization}.$$

since we're already computing the Hessian, we can use it to optimize the log posterior, with updates

$$f \leftarrow f - (\nabla \nabla^{\mathsf{T}} \log p(f_X \mid y))^{-1} \nabla \log p(f_X \mid y)$$

- ▶ Newton-Raphson optimization converges in *one step* on quadratic functions
- For non-quadratic, convex problems, it asymptotically converges with quadratic convergence rate:

$$||f_{t+1} - f^*|| \le C \cdot ||f_t - f^*||^2$$

- Newton's method has (at least on paper) no learning rate, and no other parameters to tune.
- ▶ It requires computing and decomposing the Hessian, though

# Implementing Newton Optimization

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

$$\nabla \nabla^\mathsf{T} \log p(f_X \mid \mathbf{y}) = \sum_{i=1}^n \nabla \nabla^\mathsf{T} \log \sigma(y_i f_{x_i}) - K_{XX}^{-1} \quad \text{with} \quad \frac{\partial^2 \log \sigma(y_i f_{x_i})}{\partial f_{x_0} \partial f_{x_0}} = -\delta_{ia} \delta_{ib} \underbrace{\sigma(f_{x_i})(1 - \sigma(f_{x_i}))}_{=:w_i \text{ with } 0 < w_i < 1}$$

$$=: - \operatorname{diag} \mathbf{w} - K^{-1} = -(W + K^{-1})$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i < 1}$$

$$\stackrel{\mathsf{Procedure GP-CLASSIFY-TRAIN}(K_{XX}, m_X, \mathbf{y})}{=:w_i \text{ with } 0 < w_i <$$

return f 10 end procedure



► The Newton step can be numerically unstable when  $(W + K^{-1})$  has very small eigenvalues. It is helpful to consider the matrix

$$B := I + W^{\frac{1}{2}}KW^{\frac{1}{2}}$$

which is symmetric positive definite, and has eigenvalues  $\geq 1$ .

► From the matrix inversion lemma, we get

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

and thus

$$f_{t+1} = f_t - (\nabla \nabla^{\mathsf{T}} \log p(f_t \mid y))^{-1} \nabla \log p(f_t \mid y)$$

$$= f_t + (W + K^{-1})^{-1} (\nabla \log p(y \mid f_t)) - K^{-1} f_t$$

$$= (W + K^{-1})^{-1} (\nabla \log p(y \mid f_t) + W f_t)$$

$$= (K - KW^{\frac{1}{2}}B^{-1}W^{\frac{1}{2}}K) (\nabla \log p(y \mid f_t) + W f_t)$$

From above:

$$q(f_{X} \mid y) = \int p(f_{X} \mid f_{X})q(f_{X}) df_{X} = \int \mathcal{N}(f_{X}; m_{X} + k_{XX}K_{XX}^{-1}(f_{X} - m_{X}), k_{XX} - k_{XX}K_{XX}^{-1}k_{XX})q(f_{X}) df_{X}$$

$$= \mathcal{N}(f_{X}; m_{X} + k_{XX}K_{XX}^{-1}(\hat{\mathbf{f}} - m_{X}), k_{XX} - k_{XX}K_{XX}^{-1}k_{XX} + k_{XX}K_{XX}^{-1}\hat{\Sigma}K_{XX}^{-1}k_{XX})$$

$$\log p(f_{X} \mid y)$$

with variance

$$= k_{xx} - k_{xX} K_{XX}^{-1} k_{Xx} + k_{xX} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{Xx}$$

$$= k_{xx} - k_{xX} [K_{XX}^{-1} + K_{XX}^{-1} (K^{-1} + W)^{-1} K_{XX}^{-1}] k_{Xx}$$

$$= k_{xx} - k_{xX} (K + W^{-1})^{-1} k_{xx}$$

$$= k_{xx} - k_{xX} W^{\frac{1}{2}} W^{-\frac{1}{2}} (K + W^{-1})^{-1} W^{-\frac{1}{2}} W^{\frac{1}{2}} k_{xx}$$

$$= k_{xx} - k_{xX} W^{\frac{1}{2}} B^{-1} W^{\frac{1}{2}} k_{xx}$$

ightharpoonup We can compute the pdf of  $\sigma(f)$  analytically as

$$p(\sigma(f(x))) = \mathcal{N}(\sigma^{-1}(f(x)); m_x, s_x^2) \left| \frac{\partial \sigma^{-1}(f(x))}{\partial f} \right| = \mathcal{N}(\sigma^{-1}(f); m_x, s_x^2) \frac{1}{(\sigma(1 - \sigma))}$$

- but this distribution has no analytic expected value.
- ▶ Various approximations have been proposed. A simple (but imperfect) one [MacKay, 1992] is

$$\mathbb{E}_{\mathcal{N}(f;m,s^2)}(\sigma(f)) = \sigma\left(\frac{m}{\sqrt{1 + \frac{\pi}{8}s^2}}\right)$$

▶ a more general, but also more expensive option is to build a regular "Gaussian grid" and compute an approximate integral



# Code

#### Gaussian Process Classification:

- 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

#### Please cite this course, as

```
@techreport{Tuebingen_ProbML23,
    itile =
    {Probabilistic Machine Learning},
    author = {Hennig, Philipp},
    series = {Lecture Notes
        in Machine Learning},
    year = {2023},
    institution = {Tübingen Al Center}}
```