# PROBABILISTIC MACHINE LEARNING
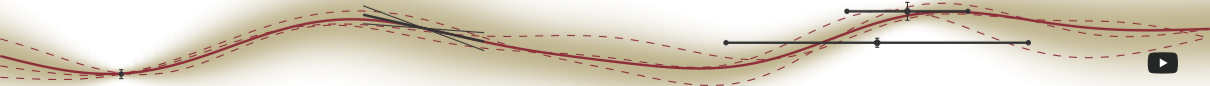## LECTURE 13
## GAUSSIAN PROCESS CLASSIFICATION

Philipp Hennig

08 June 2020

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

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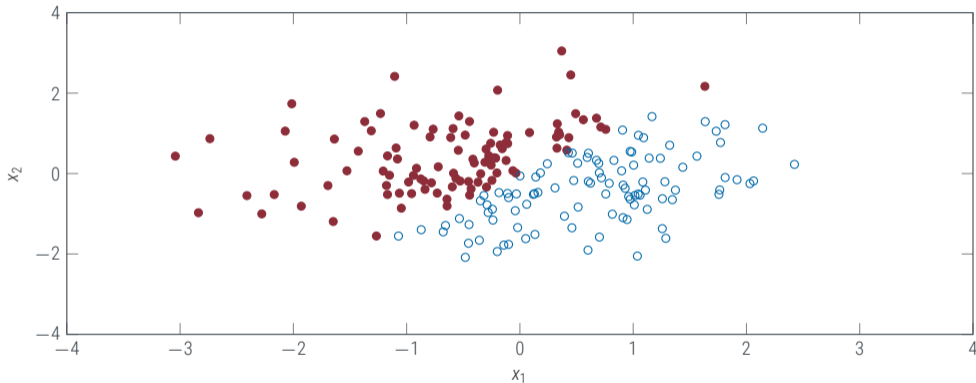CHAIR FOR THE METHODS OF MACHINE LEARNING

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 $d = 2$: binary classification)

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

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" $y_i \sim \pi_{x_i}$.

*Regression* predicts a **function**, *classification* predicts a **probability**.

Until further notice, consider only discriminative binary classification:

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

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

Until further notice, consider only discriminative binary classification:

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

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

Discriminative learning phrased probabilistically:

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

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

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

# CODE
Generative Model for Logistic Regression.ipynb

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

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

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

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

**The problem:** The posterior is not Gaussian!

$$p(f_X \mid Y) = \frac{p(Y \mid f_X) p(f_X)}{p(Y)} = \frac{\mathcal{N}(f_X; m, k) \prod_{i=1}^{n} \sigma(y_i f_{x_i})}{\int \mathcal{N}(f_X; m, k) \prod_{i=1}^{n} \sigma(y_i f_{x_i}) \, df_X}$$

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

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 $p(f, y) = p(y \mid f)p(f)$ we may be interested in

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

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

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

$Z$ for hyperparameter tuning $\qquad\qquad \bar{f}$ as a point estimator $\qquad\qquad \mathrm{var}(f)$ as an error estimator

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

Framework:

$$\int p(x_1, x_2)\, dx_2 = p(x_1) \qquad p(x_1, x_2) = p(x_1 \mid x_2)p(x_2) \qquad p(x \mid y) = \frac{p(y \mid x)p(x)}{p(y)}$$

Modelling:

▶ Directed Graphical Models
▶ Gaussian Distributions
▶ Kernels
▶ Markov Chains
▶
▶

Computation:

▶ Monte Carlo
▶ Linear algebra / Gaussian inference
▶ maximum likelihood / MAP
▶ Laplace approximations
▶

# The Laplace Approximation
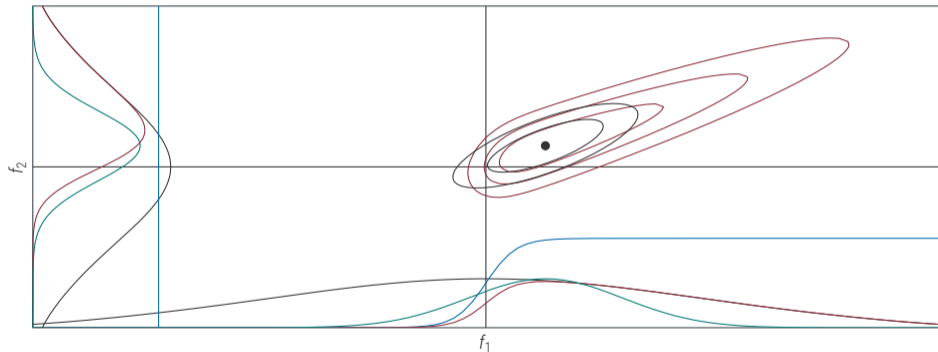
A local Gaussian approximation

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Pierre Simon M. de Laplace, 1814

# An idea as old as Probabilistic Inference
Théorie analytique des probabilités, 1814

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

▶ 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) $\log p(\theta)$

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

▶ perform **second order Taylor expansion** around $\theta = \hat{\theta} + \delta$ 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 $p$

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

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

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

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

▶ Assign approximate Gaussian posterior at training points

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

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

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

$$= \mathcal{N}(f_x;\ m_x + k_{xX}K_{XX}^{-1}(\hat{f} - m_X),\ k_{xx} - k_{xX}K_{XX}^{-1}k_{Xx} + k_{xX}K_{XX}^{-1}\hat{\Sigma}K_{XX}^{-1}k_{Xx})$$

Compare with *exact* predictions

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

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

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

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

▶ Assign approximate Gaussian posterior at training points

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

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

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

$$= \mathcal{N}(f_x;\ m_x + k_{xX}K_{XX}^{-1}(\hat{f} - m_X),\ k_{xx} - k_{xX}K_{XX}^{-1}k_{Xx} + k_{xX}K_{XX}^{-1}\hat{\Sigma}K_{XX}^{-1}k_{Xx})$$

Compare with *exact* predictions

$$\mathrm{var}_{p(f_x,f_X|y)}(f_x) = \int (f_x - \bar{f}_x)^2\,dp(f_x \mid f_X)\,dp(f_X) = k_{xx} - k_{xX}K_{XX}^{-1}k_{Xx} + k_{xX}K_{XX}^{-1}\mathrm{var}_{p(f_X|y)}(f_X)K_{XX}^{-1}k_{Xx}$$

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

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

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

▶ Assign approximate Gaussian posterior at training points

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

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

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

$$= \mathcal{N}(f_x; \ m_x + k_{xX}K_{XX}^{-1}(\hat{f} - m_X), \ k_{xx} - k_{xX}K_{XX}^{-1}k_{Xx} + k_{xX}K_{XX}^{-1}\hat{\Sigma}K_{XX}^{-1}k_{Xx})$$

▶ 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 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

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

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

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

$$\log p(f_X \mid \mathbf{y}) = \log p(\mathbf{y} \mid f_X) + \log p(f_X) - \log p(\mathbf{y}) \quad \text{with} \quad \log \sigma(y_i f_{x_i}) = -\log(1 + e^{-y_i f_{x_i}})$$

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

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

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

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

# Implementing the Laplace Approximation I
Newton Optimization

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[Rasmussen & Williams, 2006, §3.4]

```
1  procedure OPTIMIZE(L(·), f₀)
2      f ← f₀                                                    // initialize
3      while not converged do
4              g ← ∇L(f)                                         // compute gradient
5              H ← (∇∇ᵀL(f))⁻¹                                   // compute inverse Hessian
6              Δ ← Hg                                            // Newton step
7              f ← f − Δ                                         // perform step
8          converged ← ‖Δ‖ < ϵ                                  // check for convergence
9      end while
10     return f
11 end procedure
```

# Implementing the Laplace Approximation II

Concrete Algorithm, preliminary form

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[Rasmussen & Williams, 2006, §3.4]

```
1  procedure GP-LOGISTIC-TRAIN(K_XX, m_X, y)
2      f ← m_X                                                                          // initialize
3      while not converged do
4              r ← (y+1)/2 − σ(f)                          // = ∇ log p(y | f_X), gradient of log likelihood
5              W ← diag(σ(f) ⊙ (1 − σ(f)))               // = −∇∇ log p(y | f_X), Hessian of log likelihood
6              g ← r − K_XX^{-1}(f − m_X)                                             // compute gradient
7              H ← −(W + K^{-1})^{-1}                                          // compute inverse Hessian
8              Δ ← Hg                                                                 // Newton step
9              f ← f − Δ                                                             // perform step
10         converged ← ‖Δ‖ < ε                                              // 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, use $B := I + W^{1/2}K_{XX}W^{1/2}$ (cf. Rasmussen & Williams).

# Computing Predictions
from the Laplace approximation

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[Rasmussen & Williams, 2006, §3.4]

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

$$\nabla \log p(f_X \mid y) = \underbrace{\sum_{i=1}^{n} \nabla \log \sigma(y_i f_{x_i})}_{=:r} - K_{XX}^{-1}(f_X - m_X) \quad \text{with} \quad \frac{\partial \log \sigma(y_i f_{x_i})}{\partial f_{x_j}} = \delta_{ij}\left(\frac{y_i + 1}{2} - \sigma(f_{x_i})\right)$$
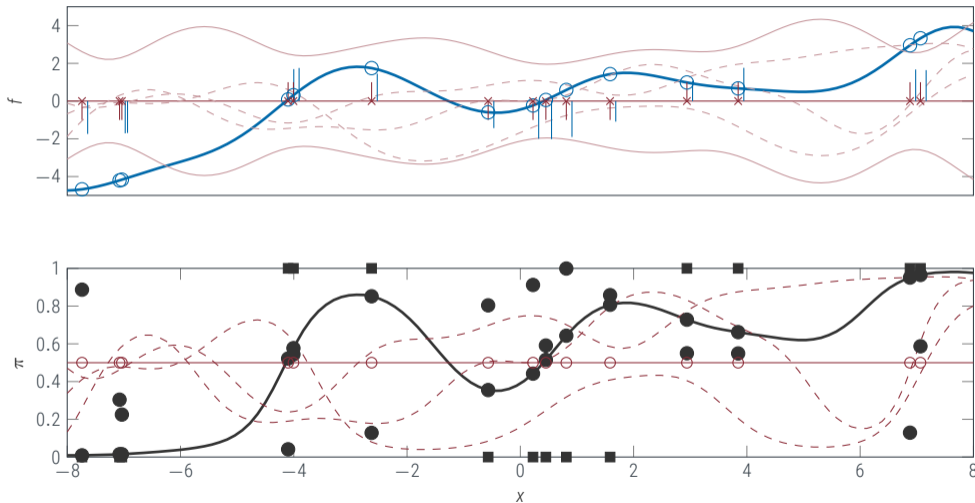
1 **procedure** GP-LOGISTIC-PREDICT($\hat{f}, W, R, r, k, x$)   $/\!/$ $\hat{f}, W, R = $ Cholesky($B$), $r$ handed over from training

2     **for** $i = 1, \ldots,$ LENGTH($x$) **do**

3        $\bar{f}_i \leftarrow k_{x_i X} r$   $/\!/$ mean prediction (note at minimum, $0 = \nabla p(f_X \mid y) = r - K_{XX}^{-1}(f_X - m_X)$).

4        $s \leftarrow R^{-1}(W^{1/2} k_{Xx_i})$   $/\!/$ pre-computation allows this step in $\mathcal{O}(n^2)$

5        $v \leftarrow k_{x_i x_i} - s^\intercal s$   $/\!/$ $v = \text{cov}(f_x)$

6        $\bar{\pi}_i \leftarrow \int \sigma(f_i) \mathcal{N}(f_i, \bar{f}_i, v)\, df_i$   $/\!/$ predictive probability for class 1 is $p(y \mid \bar{f}) = \int p(y_x \mid f_x)p(f_x \mid \bar{f})\, df_x$

7     **end for**   $/\!/$ entire loop is $\mathcal{O}(n^2 m)$ for $m$ test cases.

8     **return** $\bar{\pi}_x$

9 **end procedure**

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

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