Deadline: Fr. Mai 10 Drop your printed or legible handwritten submissions into the boxes at Samelsonplatz, or upload a Jupyter notebook (.ipynb) or a .pdf file via LearnWeb.

1 GPs & Linear Regression

(literature: Rasmussen Chapter 6, Murphy Chapter 15.4)

For a standard linear regression model the likelihood is given by

$$p(f_* \mid x_*, (X, Y)) = \mathcal{N}(Y \mid X\beta, \sigma^2 I)$$

With the optimal choice of parameters being $\hat{\beta} = (X^{\intercal}X)^{-1}X^{\intercal}Y$. A Bayesian linear regression (BLR) model (Murphy 7.6) extends this by factoring in a prior belief about β

$$p(\beta) = \mathcal{N}(\mu_{\beta}, \Sigma_{\beta})$$

yielding the following posterior for the parameters

$$p(\beta|X, y) \propto p(y|X, \beta)p(\beta)$$

= $\mathcal{N}(y \mid X\beta, \sigma^2 I)\mathcal{N}(\beta \mid \mu_{\beta}, \Sigma_{\beta}) = \mathcal{N}(\beta \mid \mu'_{\beta}, \Sigma'_{\beta})$

and subsequently the posterior predictive distribution

$$p(y_* \mid X_*) = \mathbb{E}_{\beta}[p(y \mid \beta, X_*)] = \int p(y \mid X_*, \beta) p(\beta) d\beta$$
$$= \int \mathcal{N}(y \mid X_*\beta, \sigma^2 I) \mathcal{N}(\beta \mid \mu_{\beta}', \Sigma_{\beta}') d\beta$$

A. [4p] Show that mean and variance of $p(\beta \mid X, y)$ are given by

$$\begin{split} \mu'_{\beta} &= (\frac{1}{\sigma^2} X^{\mathsf{T}} X + \Sigma_{\beta}^{-1})^{-1} (\frac{1}{\sigma^2} X^{\mathsf{T}} Y + \Sigma_{\beta}^{-1} \mu_{\beta}) \\ \Sigma'_{\beta} &= (\frac{1}{\sigma^2} X^{\mathsf{T}} X + \Sigma_{\beta}^{-1})^{-1} \end{split}$$

Note that we have already encountered such a model before! Ridge regression can be interpreted as a BLR model with $\mu_{\beta} = 0, \Sigma_{\beta} = \lambda I$.

- **B.** [4p] Show that $p(y_* \mid X_*) = \mathcal{N}(y \mid X_* \mu'_\beta, \sigma^2 I + X_* \Sigma'_\beta X_*^\intercal)$
- **C.** [4p] Show that the BLR is equivalent to a GP with a bi-linear kernel function

$$\kappa(x, x') = x^T \Lambda x'$$

Hint: You will need the following formulae

Theorem 1. If $p(x) = \mathcal{N}(x \mid \mu_x, \Sigma_x)$ and $p(y \mid x) = \mathcal{N}(y \mid Ax + b, \Sigma_y)$, then

$$p(x \mid y) = \mathcal{N}(x \mid \mu_{x\mid y}, \Sigma_{x\mid y}) \quad with \quad \begin{aligned} \mu_{x\mid y} &= \Sigma_{x\mid y} \left[A^{\mathsf{T}} \Sigma_{y}^{-1}(y-b) + \Sigma_{x}^{-1} \mu_{x} \right] \\ \Sigma_{x\mid y}^{-1} &= \Sigma_{x}^{-1} + A^{\mathsf{T}} \Sigma_{y}^{-1} A \end{aligned}$$
$$p(y) = \mathcal{N}(y \mid A\mu_{x} + b, \Sigma_{y} + A\Sigma_{x} A^{\mathsf{T}})$$

(Murphy, 4.125 & 4.126)

Theorem 2. Matrix Inversion Lemma

$$(A + UBV)^{-1} = A^{-1} - A^{-1}U (B^{-1} + VA^{-1}U)^{-1} VA^{-1}$$

2 GP Classification

- **A.** [4p] Explain briefly how classification with GPs works. What are the main difficulties?
- **B.** [4p] Explain briefly how the Laplace Approximation method works

1/**1**

(8 points)