

Rubén Martínez Cantín

Dpto. Informática e Ingeniería de Sistemas.

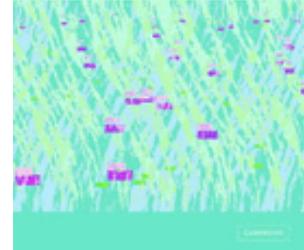
Monte Carlo

- Markov Chain Monte Carlo

- **Bayesian decision making**

- Active learning
- Bayesian optimization

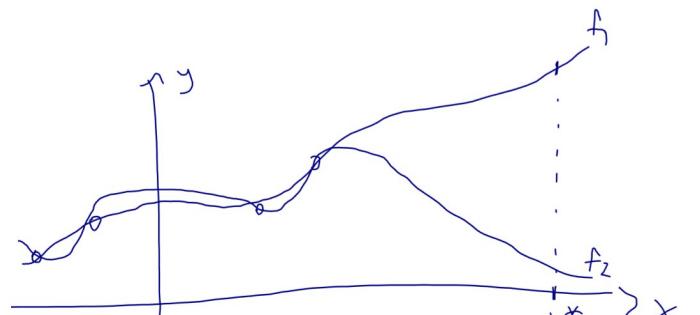
- Phillip Henning, Probabilistic Machine Learning,
<https://uni-tuebingen.de/en/180804>



- Python resources

- In the labs: GPy, GPyOpt , scikit-learn.
- Gpyflow, PyMC3...

engineer to use for your diagnosis?





"oveja"



Universidad
Zaragoza

Images: Wikipedia, Nando de Freitas.

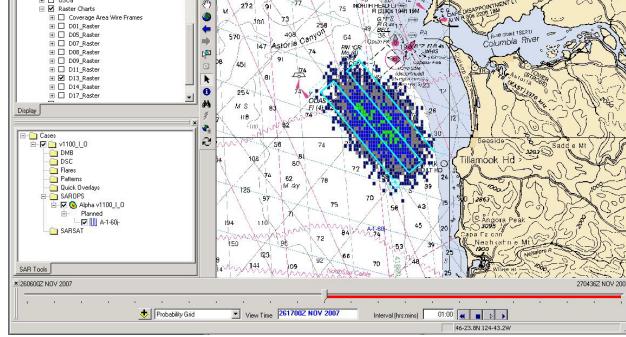
- Bayes Rule

$$P(A|B) = \frac{\int_{\mathcal{A}} P(B|A)P(A)dA}{\int_{\mathcal{A}} P(B|A)P(A)dA}$$

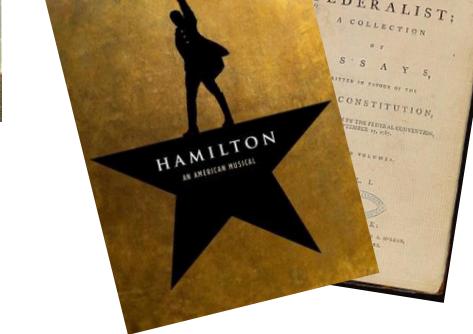
- Expectation

$$\mathbb{E}_{P(A)}f(A) = \int_{\mathcal{A}} f(A)P(A)dA$$

- Which animal do you think is deadlier (higher probability of dying if bitten)?
 - $p(\text{shark bite})?$
 - $p(\text{dead} \mid \text{shark bite})?$
 - $p(\text{shark bite} \mid \text{dead})?$



3G / 4G



Images: Wikipedia, Judea Pearl, Katie Bowman, Claude Berrou, David J.C. Mackay.

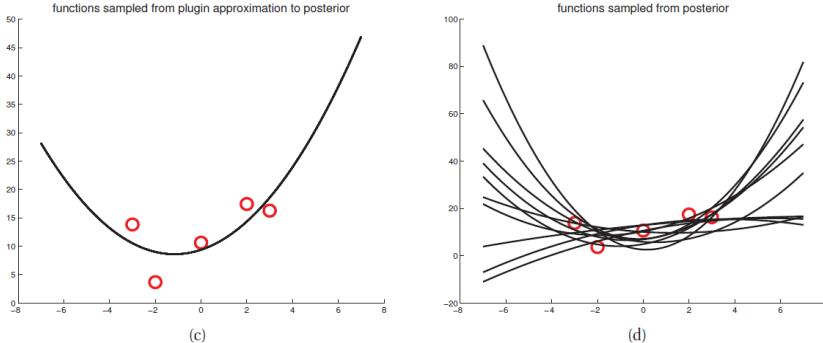


Image: Kevin P. Murphy

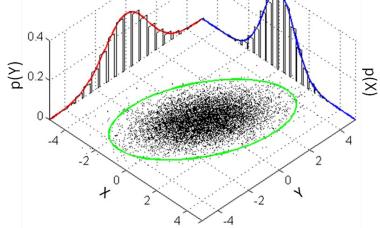
- Step 3: Compute the posterior:

$$p(\textcolor{red}{w}|\textcolor{green}{D}) = \frac{p(\textcolor{green}{D}|\textcolor{red}{w})p(\textcolor{red}{w})}{p(\textcolor{green}{D})}$$

- What is the prior?
 - For a linear-Gaussian likelihood model, the conjugate prior is a Gaussian.
 - The posterior is also Gaussian → Closed-form, easy, allows recurrence.

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{w}_0, \Sigma_0)$$

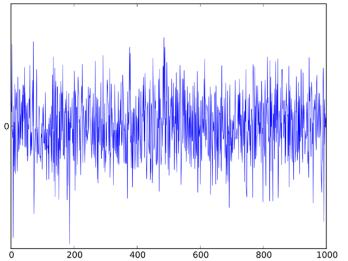
- $p(x)$ and $p(y)$ are Gaussian.
- $p(x|y)$ and $p(y|x)$ are Gaussian.
- Linear: $a * p(x) + b * p(y)$ is a Gaussian.
- Product: $p(x)*p(y)$ is a Gaussian:
- Convolution: $p(x + y)$ is a Gaussian



Images: Wikipedia

$$\begin{aligned} p(\mathbf{x}_1) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \\ p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{x}_1 | \mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}) \\ \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\Sigma}_{1|2} (\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1} \end{aligned}$$



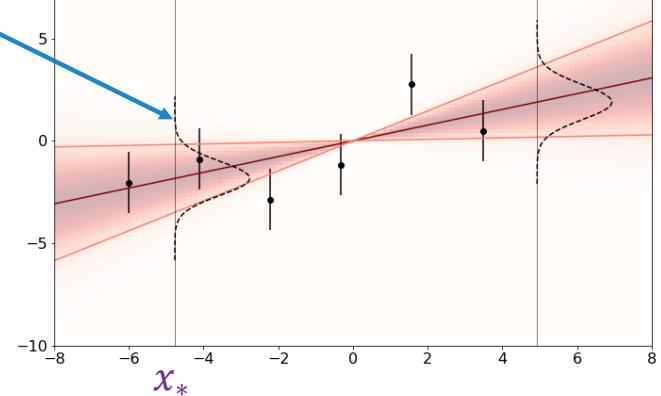
- Noise is usually Gaussian

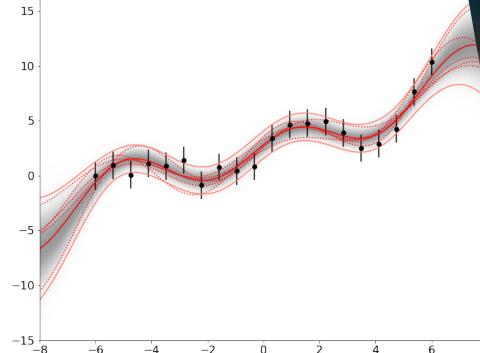
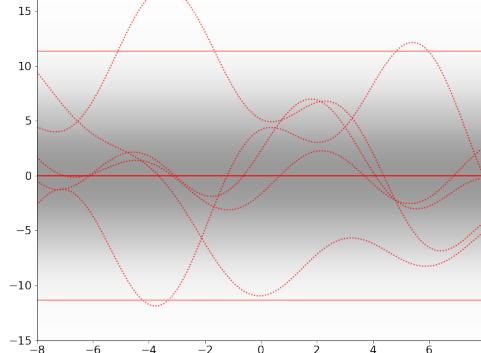
- If we use an isotropic Gaussian prior: $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{0}, \tau^2 I_d)$ then

$$y = \mathbf{x}^T \mathbf{w} + \epsilon$$

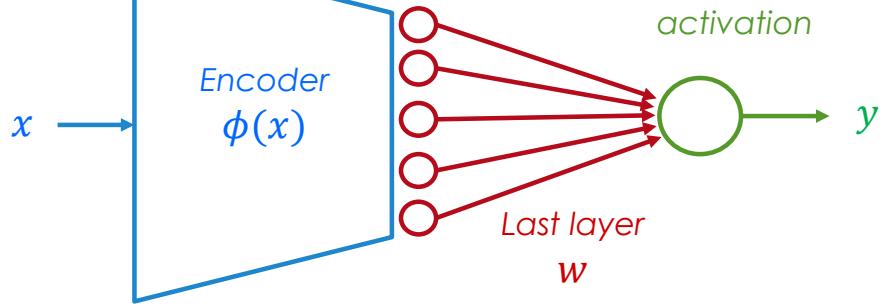
$$\mu(\mathbf{x}_*) = ?$$

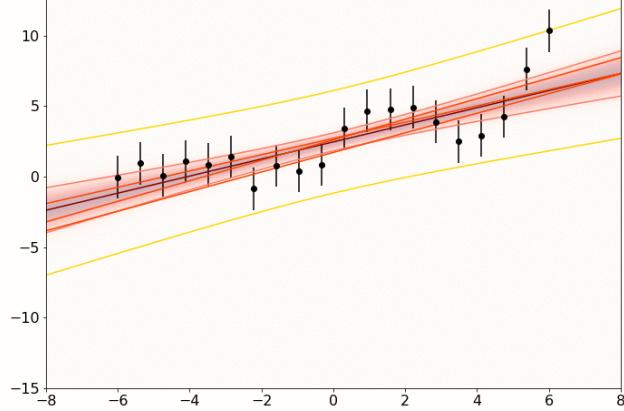
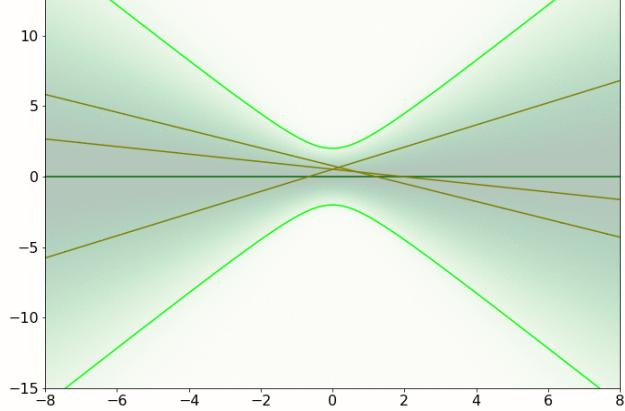
$$\sigma_n^2(\mathbf{x}_*) = ?$$

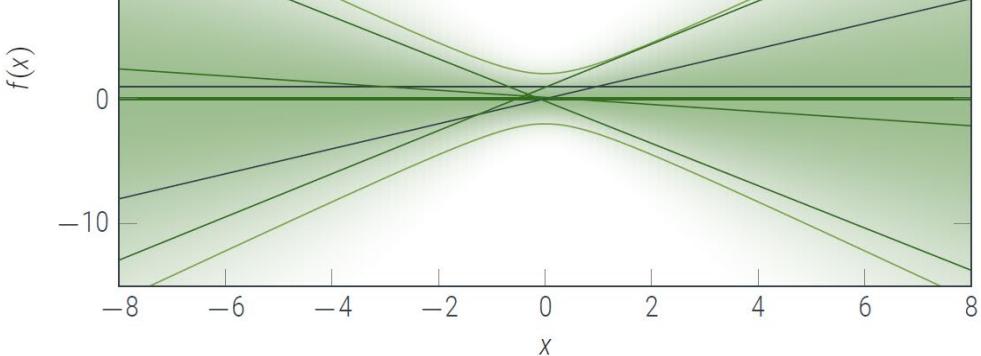
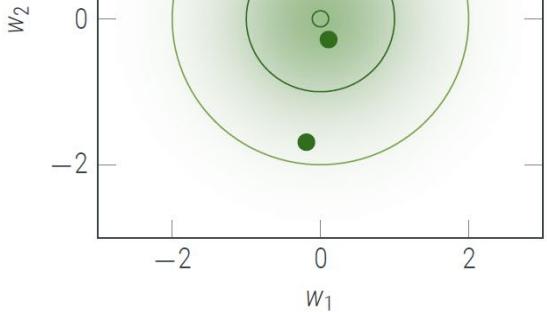


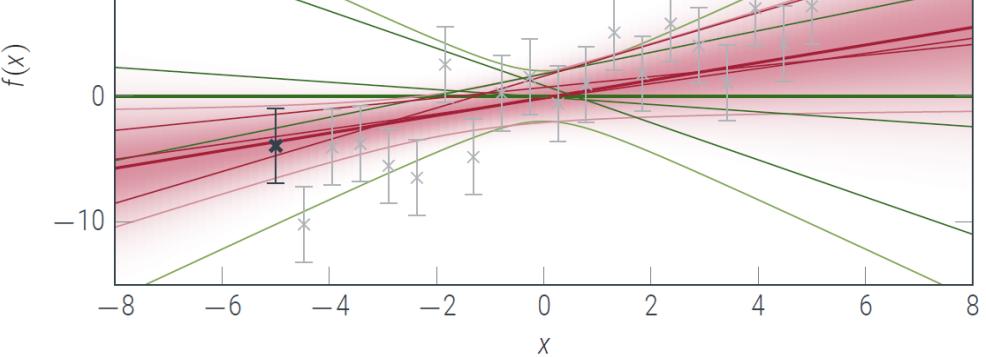
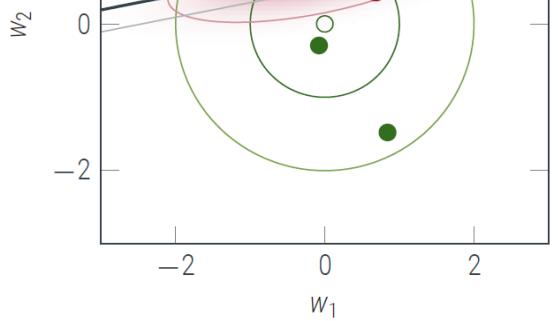


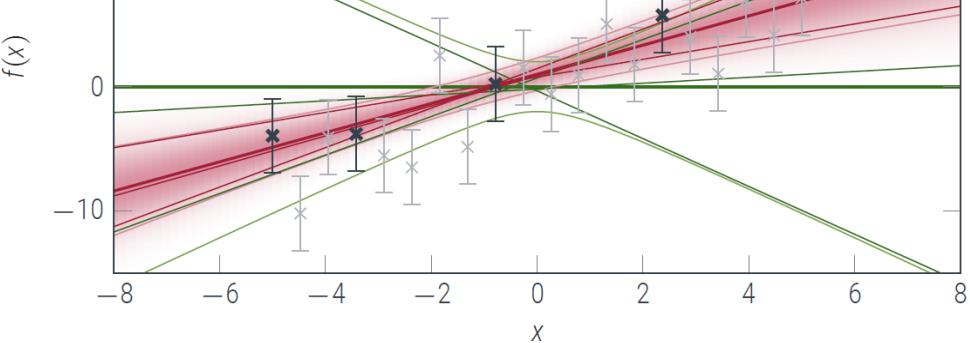
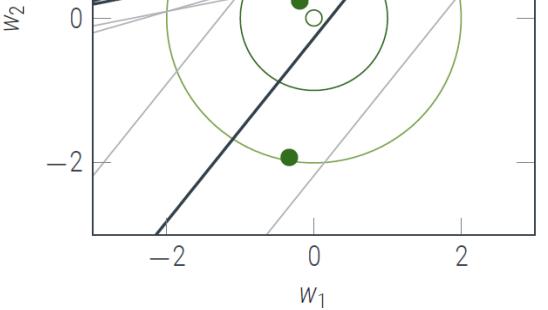
Universidad
Zaragoza

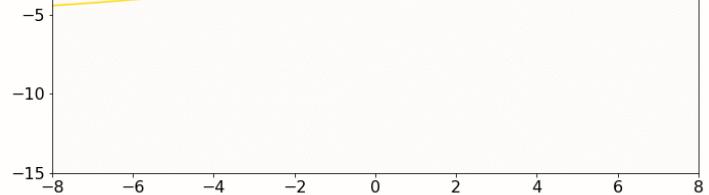
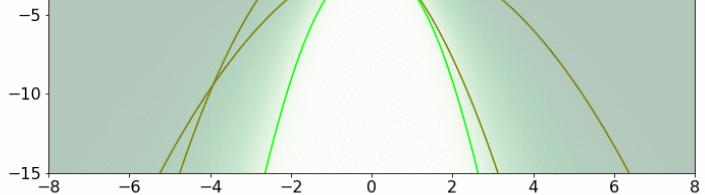


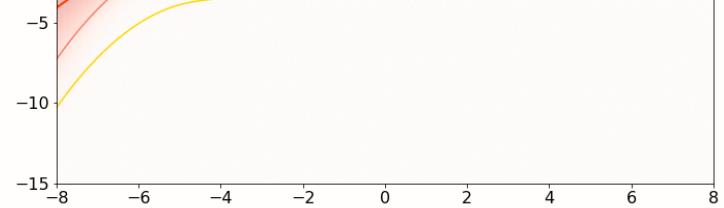
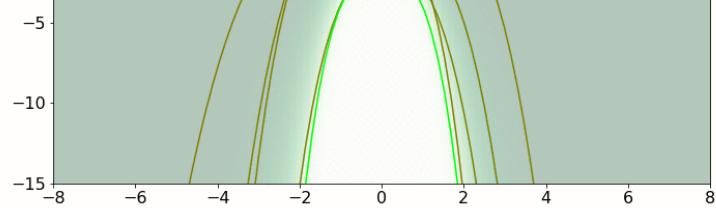


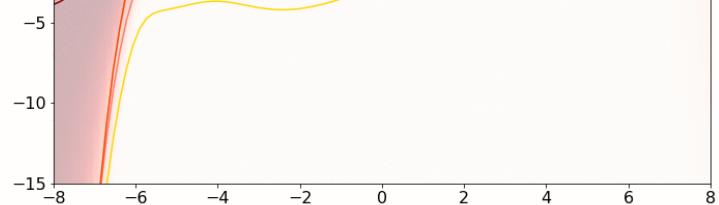
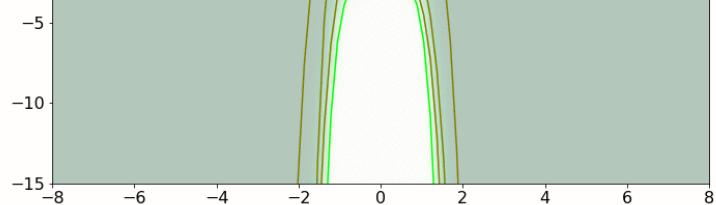


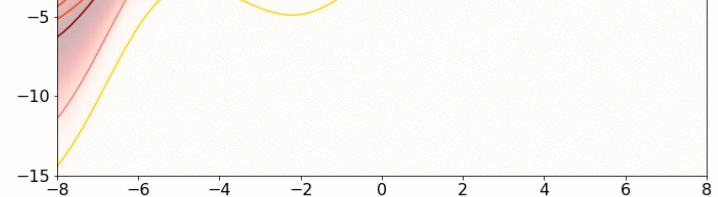
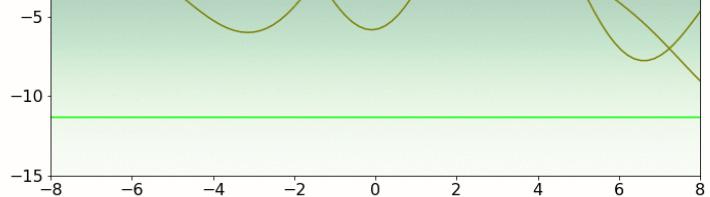


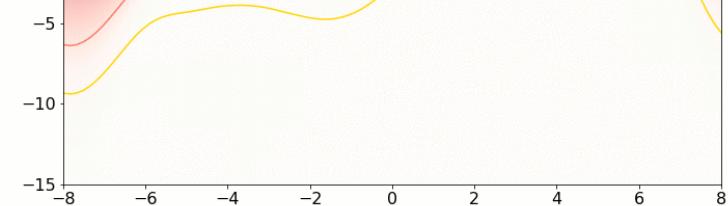
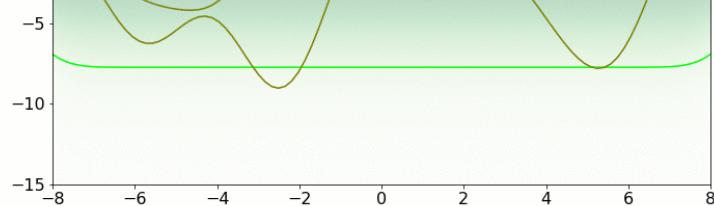


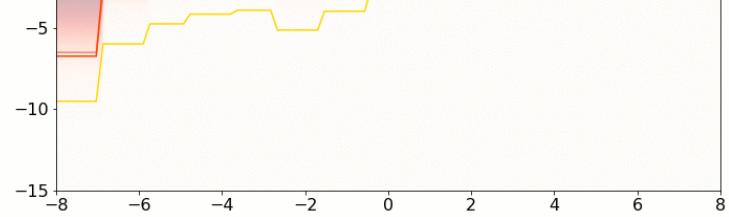
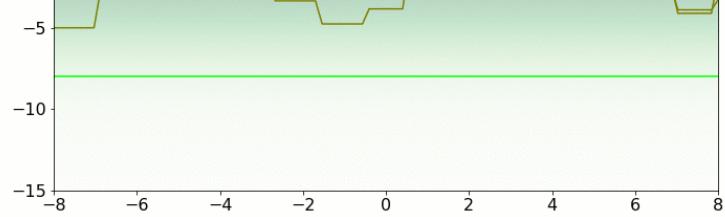


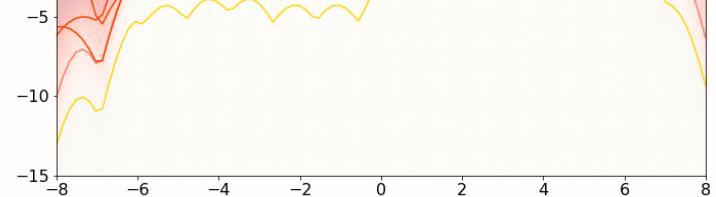
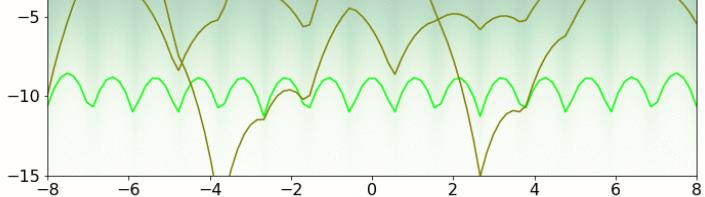












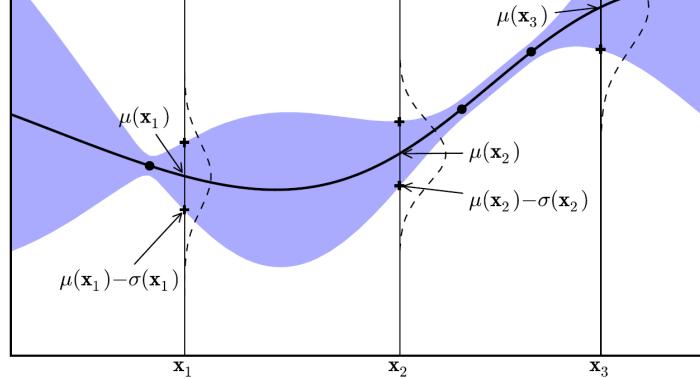
- Hyperparameters are nonlinear.
 - No closed-form. No Gaussian...
- We'll get back to it latter.

result is always a scalar.

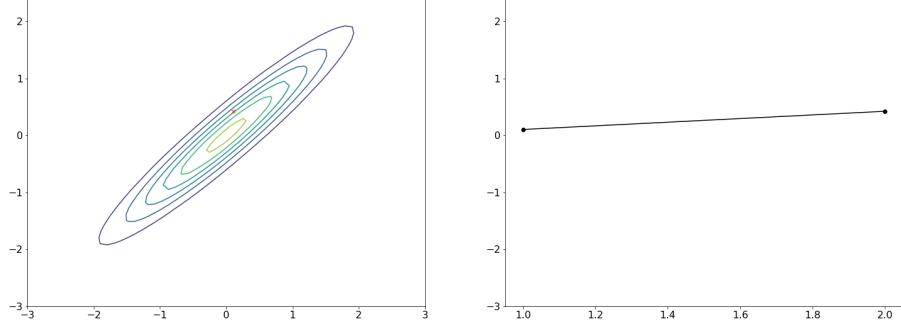
- Kernel function: $k(x_1, x_2) = \phi(x_1)^T \Sigma \phi(x_2)$
 - This type of function is called a Mercer kernel.

- The predictive posterior is a conditional Gaussian of:

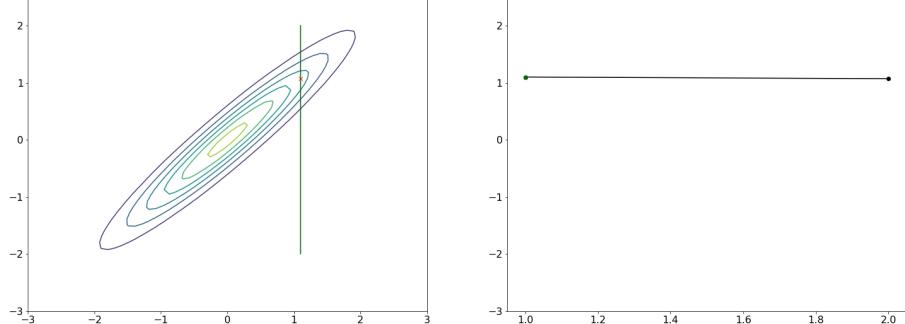
$$\begin{pmatrix} \mathbf{y} \\ f_* \end{pmatrix} \sim N \left(0, \begin{bmatrix} k(X, X) + \sigma^2 I & k(X, x_*) \\ k(x_*, X) & k(x_*, x_*) \end{bmatrix} \right)$$



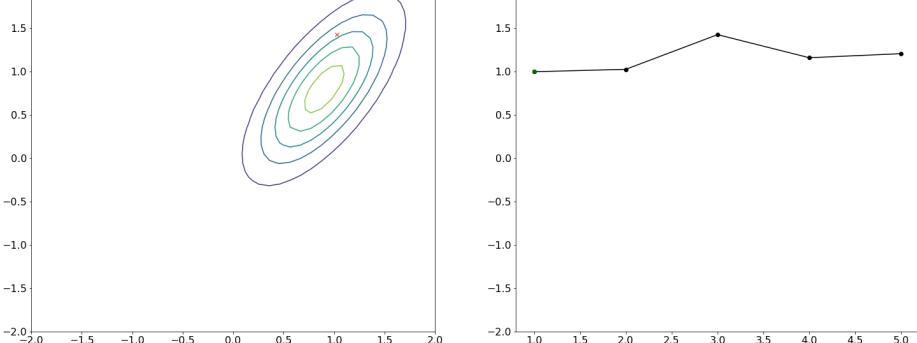
Large correlation

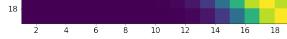


Large correlation

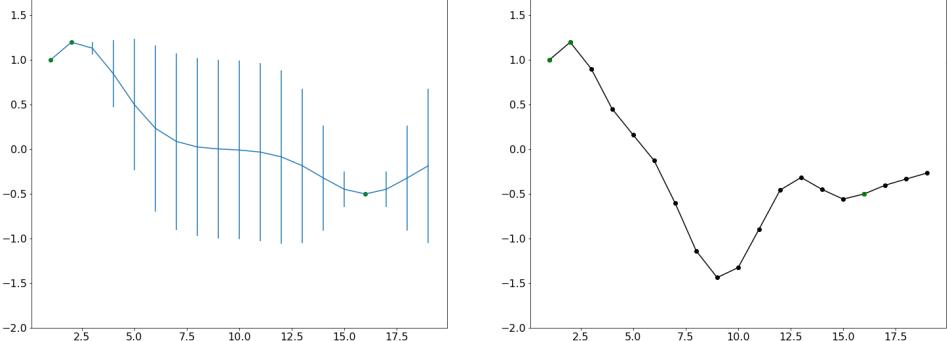


Conditional





Conditional



data points.

- GPs are distributions over functions:

$$f(x) \sim \mathcal{GP}(m(x), k(x, x'))$$

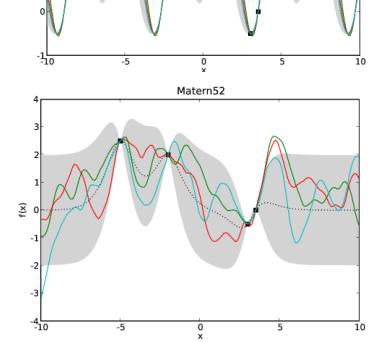
$$m(x) = \mathbb{E}[f(x)]$$

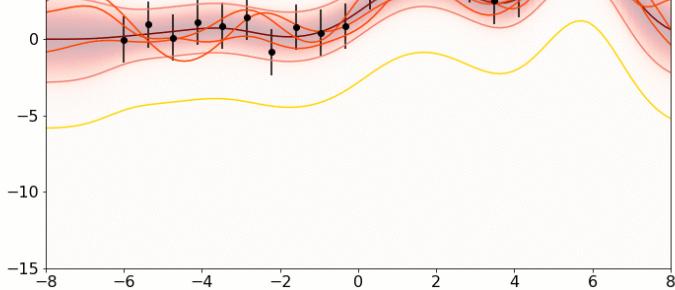
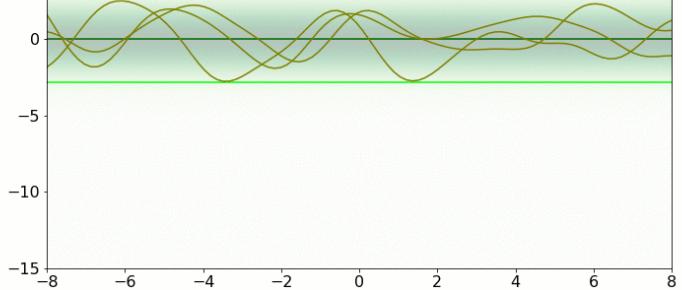
$$k(x, x') = \mathbb{E} \left[(f(x) - m(x))(f(x) - m(x))^T \right]$$

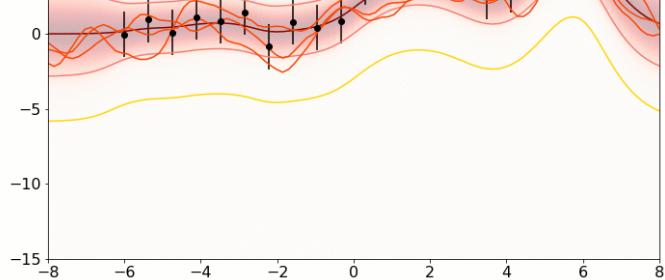
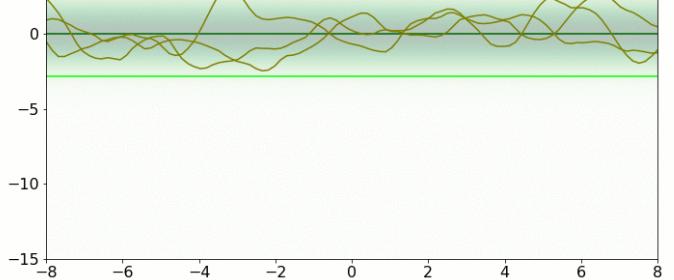
rational quadratic
neural network

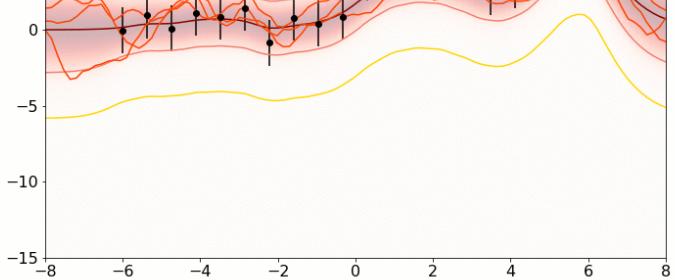
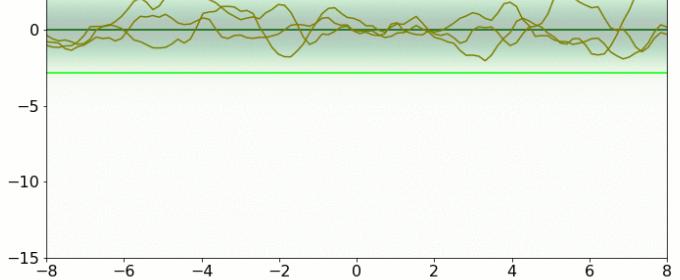
$$(1 + \frac{r^2}{2\alpha\ell^2})^{-\alpha}$$
$$\sin^{-1} \left(\frac{2\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}}'}{\sqrt{(1+2\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}})(1+2\tilde{\mathbf{x}}'^\top \Sigma \tilde{\mathbf{x}}')}} \right)$$

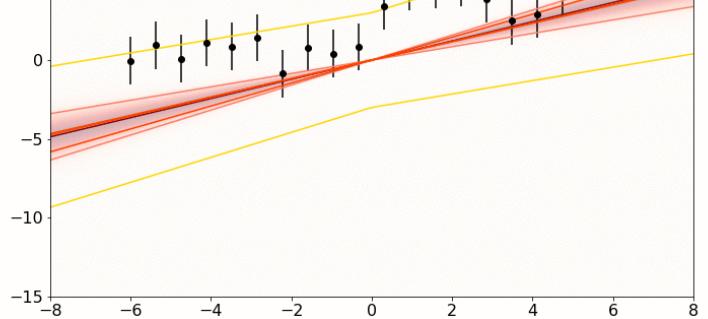
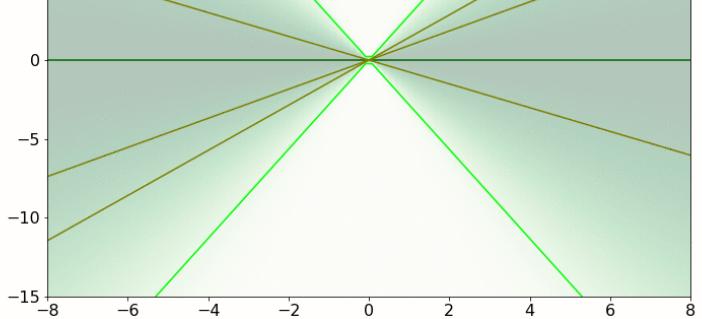
- Others: Wiener process, splines...
- Not only real vectors: strings, paths...

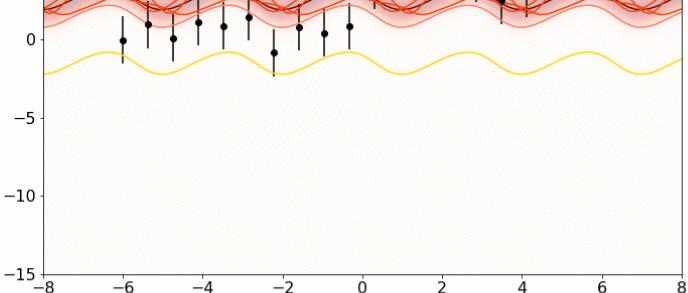
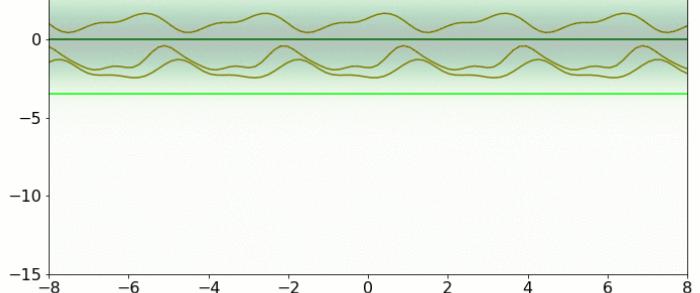


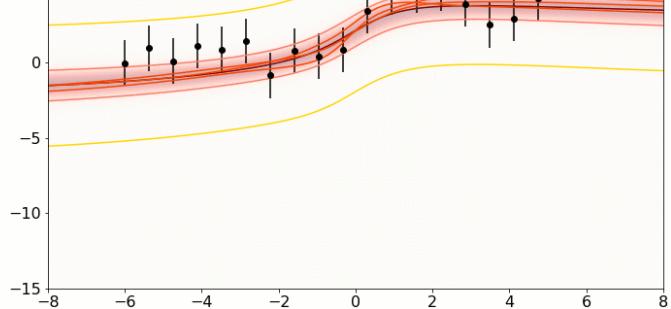
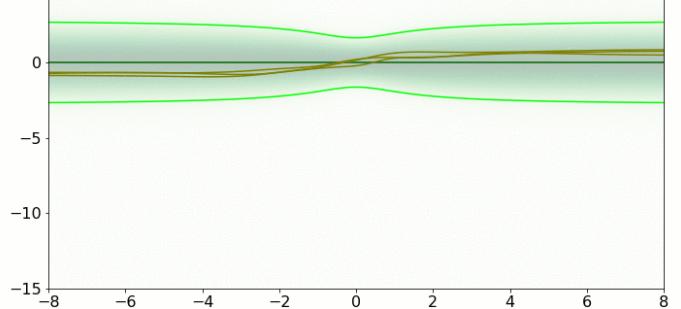




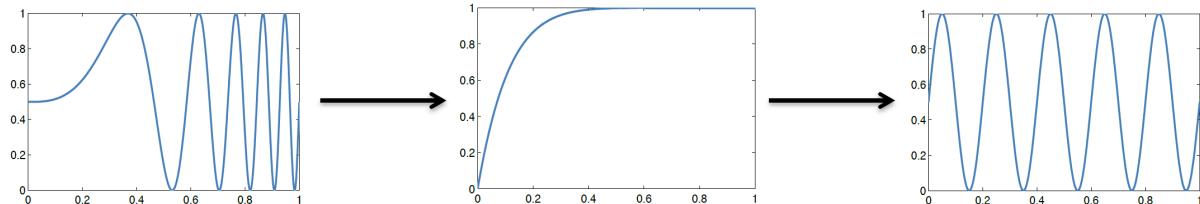


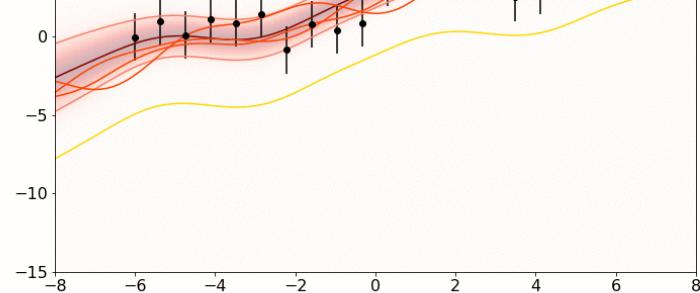
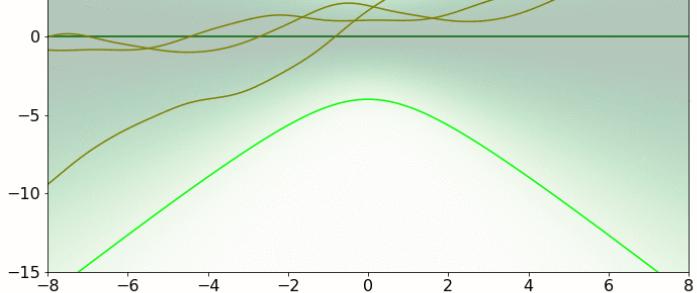


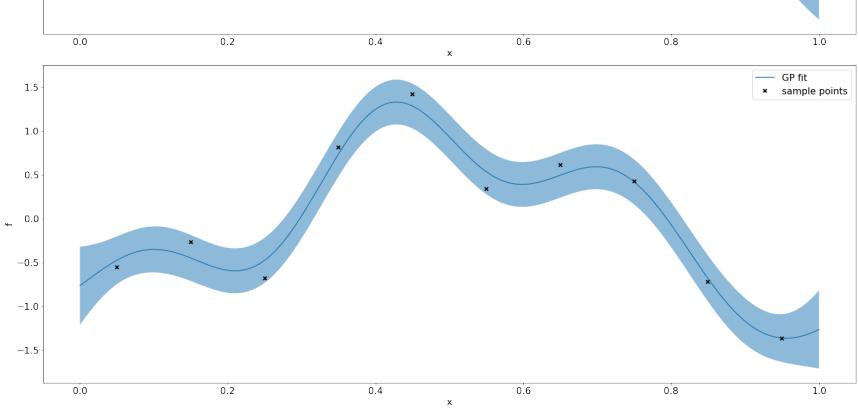




- Mapped inputs: $k(\psi(x), \psi(x'))$



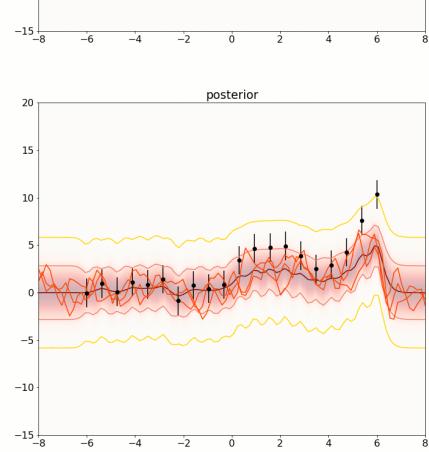
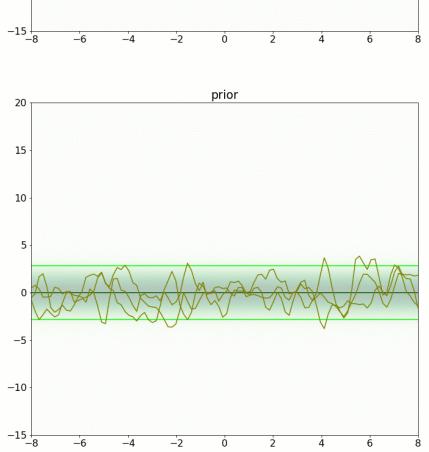




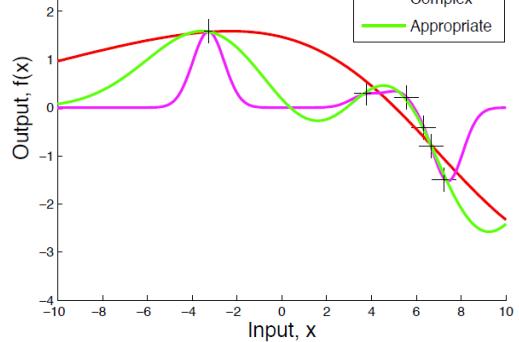
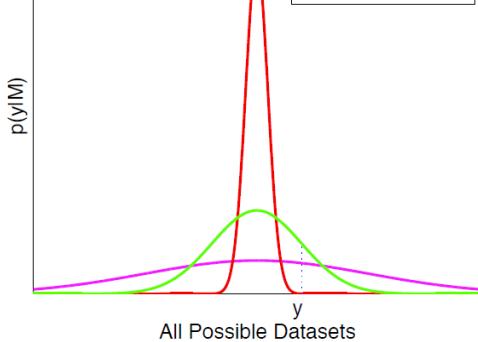


Universidad
Zaragoza

$\ell=0.3$



- Sampling (e.g.: MCMC – next lesson):
 - *Full Bayesian*
 - Result in a mixture of GPs
 - Much, much more expensive



- Learn – Maximize the marginal log-likelihood:

$$\log p(\mathbf{y} | X, \theta) = -\frac{1}{2}\mathbf{y}^T \mathbf{K}_\theta^{-1} \mathbf{y} - \frac{1}{2}\log|\mathbf{K}_\theta| - \frac{n}{2}\log 2\pi$$

\mathbf{x}_* (test input)

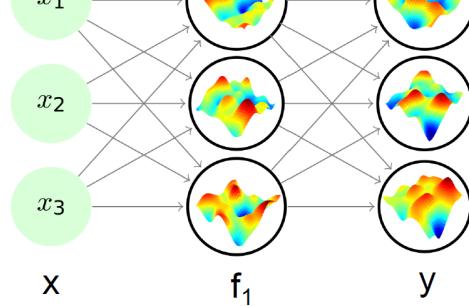
- 2: $L := \text{cholesky}(K + \sigma_n^2 I)$
 - 3: $\boldsymbol{\alpha} := L^\top \backslash (L \backslash \mathbf{y})$
 - 4: $\bar{f}_* := \mathbf{k}_*^\top \boldsymbol{\alpha}$
 - 5: $\mathbf{v} := L \backslash \mathbf{k}_*$
 - 6: $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$
 - 7: $\log p(\mathbf{y}|X) := -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{n}{2} \log 2\pi$
 - 8: **return:** \bar{f}_* (mean), $\mathbb{V}[f_*]$ (variance), $\log p(\mathbf{y}|X)$ (log marginal likelihood)
- } predictive mean eq. (2.25)
} predictive variance eq. (2.26)
eq. (2.30)

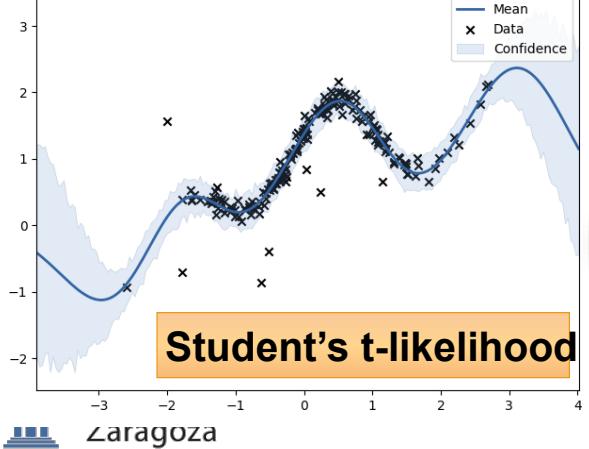


- $\alpha = (\mathbf{L} \setminus \mathbf{k}_*)^T$
- $\mu = \alpha \mathbf{y}$
- $\sigma_n^2 = k(x_*, x_*) - \alpha \mathbf{k}_*^T$
- $\beta = (\mathbf{L} \setminus \mathbf{y})$
- $\log(p(y|X)) = -\frac{1}{2}\beta^T\beta - \sum_{i=1}^N \log(L_{ii}) - \frac{N}{2}\log(2\pi)$

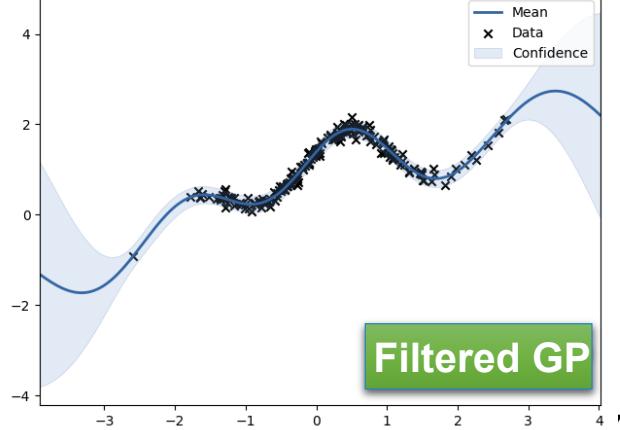
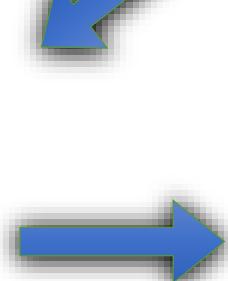
course

- Deep GPs
 - Build a model made of layers of GPs
- Robust regression:
 - Use non-Gaussian likelihoods (e.g.: Student-t)





Zaragoza
1542



- GPs for decision making
 - GPs in reinforcement learning/bandits
 - GPs for dynamic systems
 - GPs for experimental design and optimization → Meta-learning



Universidad
Zaragoza

1542

Thesis presented for the
Degree of M.Sc. in Engineering,
University of the Witwatersrand

JOHANNESBURG,
15th March, 1951.

