

# ***ADVANCED MACHINE LEARNING***

## **Non-linear regression techniques Part – II**

### **Gaussian Process Regression**



# Probabilistic Regression (PR)

## Creating the model

PR is a statistical approach to classical linear regression that estimates the relationship between zero-mean variables  $y$  and  $x$  by building a linear model:

$$y = f(x, w) = w^T x, \quad w, x \in \mathbb{R}^N$$



# Probabilistic Regression (PR)

## Creating the model

If one assumes that the observed values of  $y$  differ from  $f(x)$  by an additive noise  $\varepsilon$  that follows a zero-mean Gaussian distribution (such an assumption consists of putting a *prior distribution* over the noise), then:

$$y = w^T x + \varepsilon, \quad \text{with } \varepsilon = N(0, \sigma_\varepsilon^2)$$

The addition of a Normal random variable to a constant variable leads  $y$  to become a Normal random variable

Where have we seen this before?

Answer: RVM / RVR



# Probabilistic Linear Regression

## Maximum Likelihood Estimation

Training set of  $M$  pairs of data points  $\{X, \mathbf{y}\} = \{x^i, y^i\}_{i=1}^M$ .

Likelihood of the regressive model

$$\mathbf{y} = \underline{w^T X} + \underline{N(0, \sigma_\varepsilon^2)}$$

$$\Rightarrow \mathbf{y} \sim \underline{p(\mathbf{y} | X, \underline{w}, \underline{\sigma_\varepsilon})}$$

Parameters of  
the model.

Assume that the data points are independently and identically distributed (i.i.d), the likelihood is:

$$\begin{aligned} p(\mathbf{y} | X, w, \sigma_\varepsilon) &\sim \prod_{i=1}^M p(y^i | x^i, w, \sigma_\varepsilon) \\ &= \prod_{i=1}^M \frac{1}{\sqrt{2\pi}\sigma_\varepsilon} \exp\left(-\frac{(y^i - w^T x^i)^2}{2\sigma_\varepsilon^2}\right) \end{aligned}$$



# Probabilistic Linear Regression

## Maximum Likelihood Estimation

Training set of  $M$  pairs of data points  $\{X, \mathbf{y}\} = \{x^i, y^i\}_{i=1}^M$ .

Likelihood of the regressive model

$$\mathbf{y} = \mathbf{w}^T X + N(0, \sigma_\varepsilon^2)$$

$$\Rightarrow \mathbf{y} \sim p(\mathbf{y} | X, \mathbf{w}, \sigma_\varepsilon)$$

Fix  $\sigma_\varepsilon$

$$\mathbf{w}_{\text{MLE}} = \arg \max_{\mathbf{w}} p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_\varepsilon^2 \mathbf{I}_M)$$

Closed-form solution:

$$\mathbf{w}_{\text{MLE}} = (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{y} \rightarrow \text{Same with the OLS estimator}$$

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{N} \left( \mathbf{w}_{\text{MLE}}^T \mathbf{X} - \mathbf{y} \right)^T \left( \mathbf{w}_{\text{MLE}}^T \mathbf{X} - \mathbf{y} \right) \rightarrow \text{Mean squared prediction error}$$



# Probabilistic Linear Regression

## Maximum a Posteriori (MAP) Estimation

Training set of  $M$  pairs of data points  $\{X, \mathbf{y}\} = \{x^i, y^i\}_{i=1}^M$ .

Likelihood of the regressive model

$$\mathbf{y} = w^T X + N(0, \sigma_\varepsilon^2)$$

$$\Rightarrow \mathbf{y} \sim p(\mathbf{y} | X, w, \sigma_\varepsilon)$$

Hyperparameter  
set by the user.

$w$  is treated as a random variable.  
Set a prior on distribution of  $w$ :

$$p(w) = N(0, \Sigma_w) \propto \exp\left(-\frac{1}{2} w^T \Sigma_w^{-1} w\right)$$



# Probabilistic Linear Regression

## Maximum a Posteriori (MAP) Estimation

Prior distribution over the weights:

$$\mathbf{w} \sim \mathcal{N}(\mu_{\mathbf{w}}, \Sigma_{\mathbf{w}})$$

Model Likelihood:

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_{\epsilon}^2 \mathbf{I}_M)$$

Bayes Rule:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$$

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_{\epsilon}^2) p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})}$$

Prior and likelihood are conjugate distributions. The posterior has a closed form solution:

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\mu_{\mathbf{w}|\mathbf{X},\mathbf{y}}, \mathbf{A}^{-1})$$



# Probabilistic Linear Regression

## Maximum a Posteriori (MAP) Estimation

MAP estimates derive from the expected value of the posterior distribution:

$$\begin{aligned}\mathbf{w}_{\text{MAP}} &= \mathbb{E} \left\{ p(\mathbf{w} \mid \mathbf{x}, \mathbf{Y}) \right\} \\ \Sigma_{\text{MAP}}^{-1} &= \mathbf{A} = \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}\mathbf{X}^T + \Sigma_{\mathbf{w}}^{-1} \\ \mathbf{w}_{\text{MAP}} &= \mu_{\mathbf{w}|\mathbf{X},\mathbf{y}} = \mathbf{A}^{-1} \Sigma_{\mathbf{w}}^{-1} \mu_{\mathbf{w}} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A}^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

In the special case when the distribution is zero mean and  $\Sigma_{\mathbf{w}} = \tau \mathbf{I}$ , probabilistic regression reduces to ridge regression with  $\lambda = \frac{\sigma_{\epsilon}^2}{\tau}$ .

Ridge optimal regressor:  $\mathbf{w}^* = \left( \mathbf{X}\mathbf{X}^T + \lambda \mathbf{I} \right)^{-1} \mathbf{X}\mathbf{y}$





# Probabilistic Linear Regression

## Posterior Predictive Distribution

The predictive distribution is a distribution over output  $y$ .

$$p(\underbrace{y}_{\text{output}} \mid \underbrace{x, X, \mathbf{y}}_{\text{data}}) = \int \underbrace{p(y \mid x, w)}_{\text{Likelihood}} \underbrace{p(w \mid X, \mathbf{y})}_{\text{Posterior}} dw$$

The integral has a closed form solution in the case of Normal/Gauss distributions

$$\underline{p(y \mid x, X, \mathbf{y}) \sim N(\mu_y(x), \sigma_y^2(x))},$$

$$\mu_y(x) = x^T \underline{\mathbf{w}_{\text{MAP}}} = \frac{1}{\sigma^2} x^T A^{-1} X \mathbf{y},$$

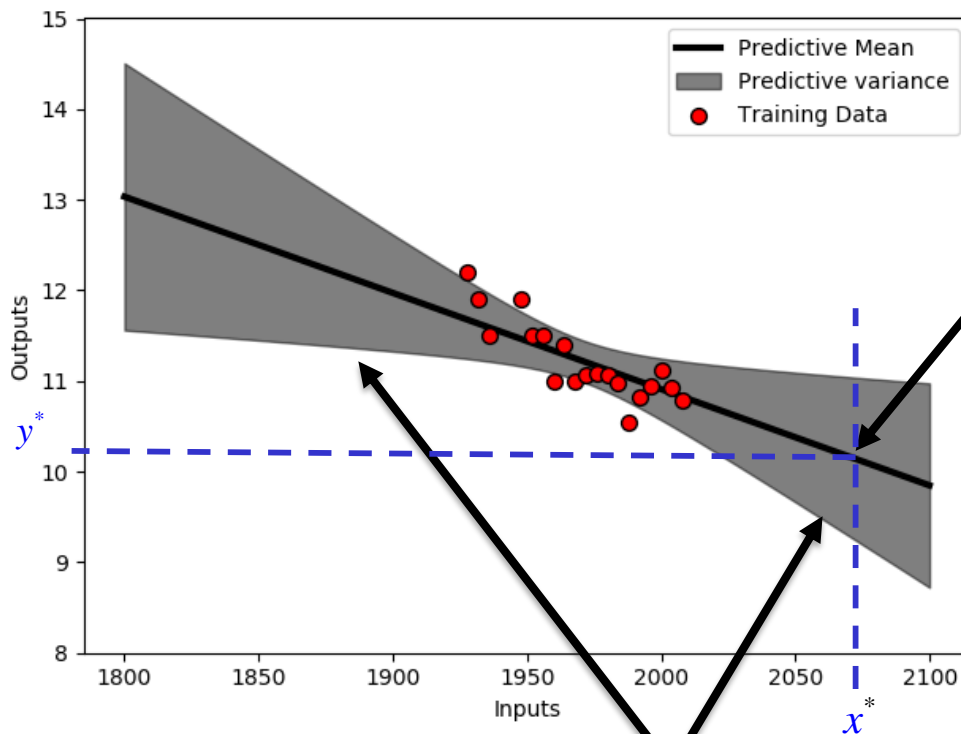
$$\sigma_y^2(x) = x^T A^{-1} x$$

$$A = \frac{1}{\sigma^2} X X^T + \Sigma_w^{-1}$$



# Probabilistic Linear Regression

## Posterior Predictive Distribution



Estimate  $y$  given a test point  $x^*$  :

$$y^* = E\{p(y | x^*, X, \mathbf{y})\} = \frac{1}{\sigma^2} x^{*T} A^{-1} X \mathbf{y}$$

Training datapoints

Testing point

The variance gives a measure of the uncertainty of the prediction:

$$\text{var}\{p(y | x)\} = x^T A^{-1} x$$



# Probabilistic Linear Regression

## Marginal Likelihood

Bayesian regression depends on hyperparameter for priors on noise and  $w$ .

The marginal likelihood provides a metric to evaluate how well our model explains the observed data and can be used to select these hyperparameters.

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_\epsilon^2) p(\mathbf{w})}{\underbrace{p(\mathbf{y} | \mathbf{X})}_{\text{Marginal Likelihood}}}$$

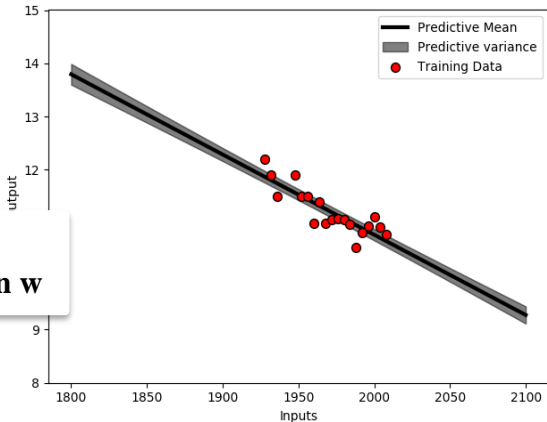
$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_\epsilon^2) p(\mathbf{w}) d\mathbf{w}$$

With a Gaussian prior and Gaussian likelihood, the marginal is also Gaussian and has a closed form solution.



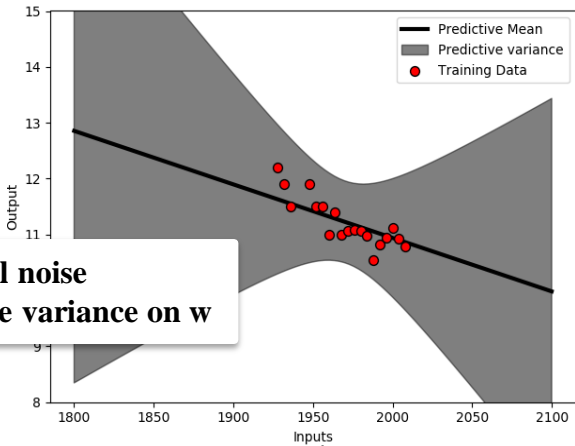
# Probabilistic Linear Regression

## Marginal Likelihood



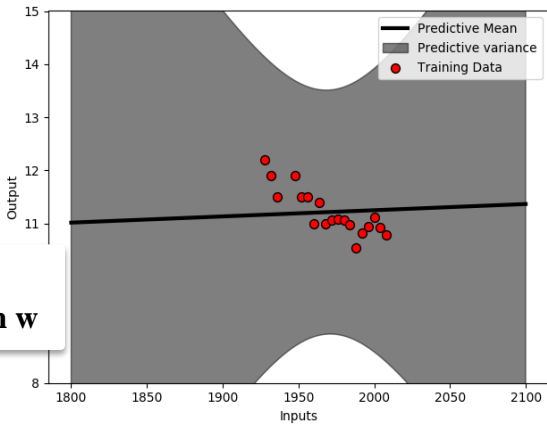
Small noise  
Small variance on  $w$

a)  
33.38



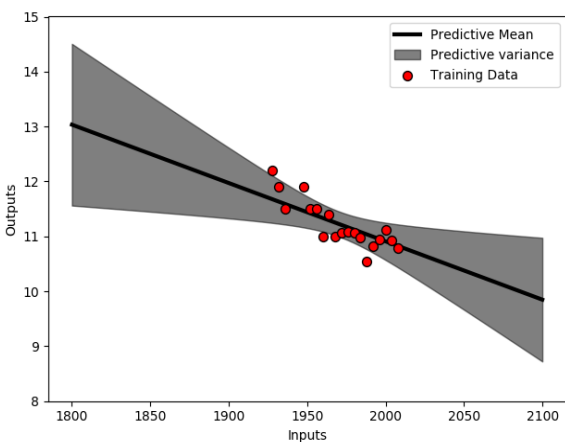
Small noise  
Large variance on  $w$

b)  
47.19



Large noise  
Large variance on  $w$

c)  
67.19



c)  
27.32

Negative Log-Marginal Likelihood for different models



# Gaussian Process Regression (GPR)

## From linear to nonlinear probabilistic regression

Linear Probabilistic Regression

$$y = \underline{w^T x} + N(0, \sigma_\varepsilon^2)$$



$$p(y | x, X, \mathbf{y}) = N\left(\frac{1}{\sigma_\varepsilon^2} \underline{x^T A^{-1} X \mathbf{y}}, x^T A^{-1} x\right),$$

$$A = \frac{1}{\sigma_\varepsilon^2} X X^T + \Sigma_w^{-1}$$

Linear Probabilistic Regression in feature space

$$y = w^T \underline{\phi(x)} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_\varepsilon^2)$$



$$p(y | x, X, \mathbf{y}) = N\left(\frac{1}{\sigma_\varepsilon^2} \underline{\phi(x)^T A^{-1} \Phi(X) \mathbf{y}}, \phi(x)^T A^{-1} \phi(x)\right)$$

with  $A = \sigma_\varepsilon^{-2} \Phi(X) \Phi(X)^T + \Sigma_w^{-1}$

**Inner product in feature space**



# Gaussian Process Regression (GPR)

## From linear to nonlinear probabilistic regression

Define the kernel as:  $k(x, x') = \phi(x)^T \Sigma_w \phi(x')$

and apply on the mean and the variance

$$y = E\{p(y | x, X, \mathbf{y})\} = \sum_{i=1}^M \alpha_i k(x, x^i)$$

$$\text{with } \alpha = [K(X, X) + \sigma_\varepsilon^2 I]^{-1} \mathbf{y}$$

$$\alpha_i > 0 \quad \forall i$$

→ All datapoints are used in the computation!

$$\text{var}(p(y | x)) = k(x, x) - K(x, X) [K(X, X) + \sigma_\varepsilon^2 I]^{-1} K(X, x)$$



# Gaussian Process Regression (GPR)

## Hyperparameters

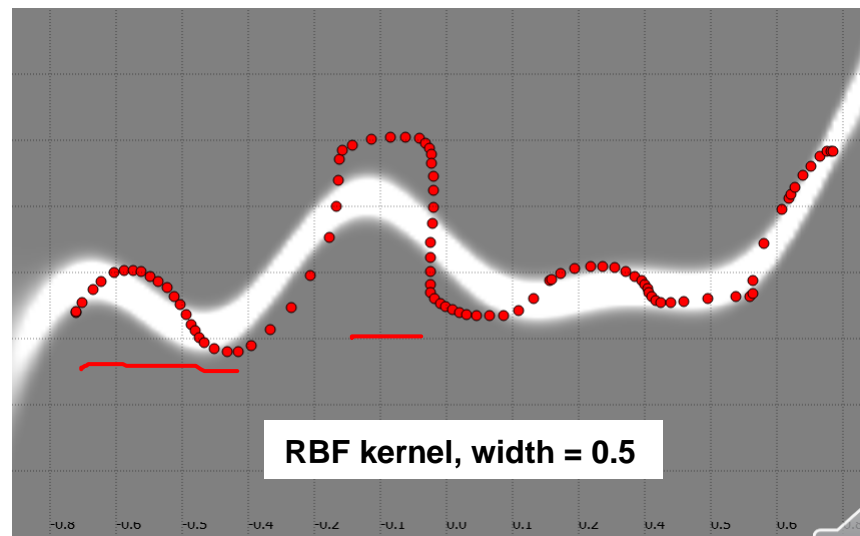
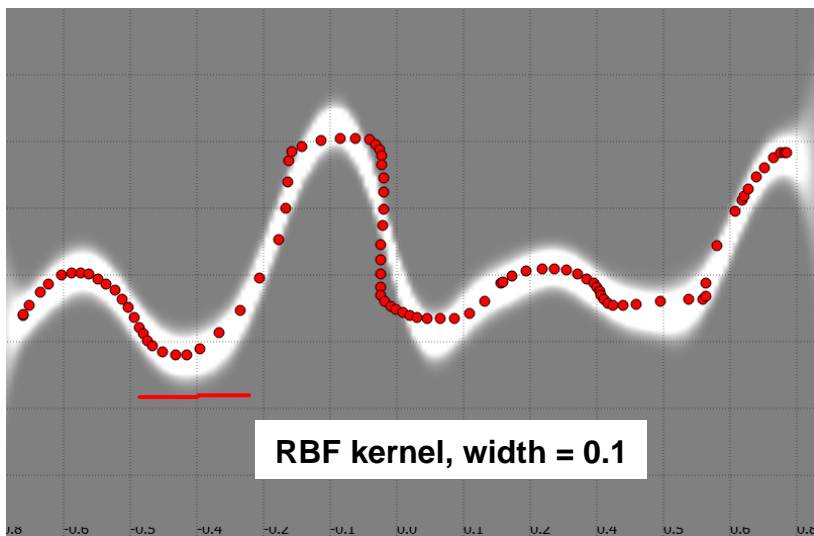
The choice of kernel and its hyperparameters will strongly influence the goodness of the fit.

$$y = \sum_{i=1}^M \alpha_i \underline{k(x, x^i)}$$

$$\text{with } \alpha = \left[ \underline{K(X, X)} + \sigma_\varepsilon^2 I \right]^{-1} y$$

$$k(x, x') = e^{-\frac{\|x-x'\|^2}{2l^2}}$$

$l$ : lengthscale parameter  
~ Kernel Width



# Gaussian Process Regression (GPR)

## Hyperparameters

The value for the **noise** needs to be pre-set by hand.  
It influence estimate of the expectation and variance of the model.

$$y = \sum_{i=1}^M \alpha_i k(x, x^i)$$
$$\text{with } \alpha = \left[ K(X, X) + \underline{\sigma_\epsilon}^2 I \right]^{-1} \mathbf{y}$$

$$\text{var}(p(y|x)) = K(x, x) - K(x, X) \left[ K(X, X) + \underline{\sigma_\epsilon}^2 I \right]^{-1} K(X, x)$$

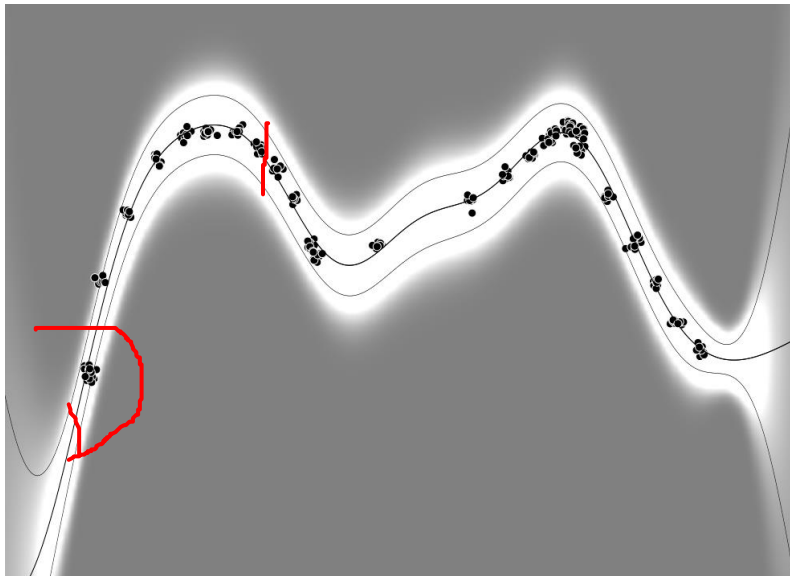




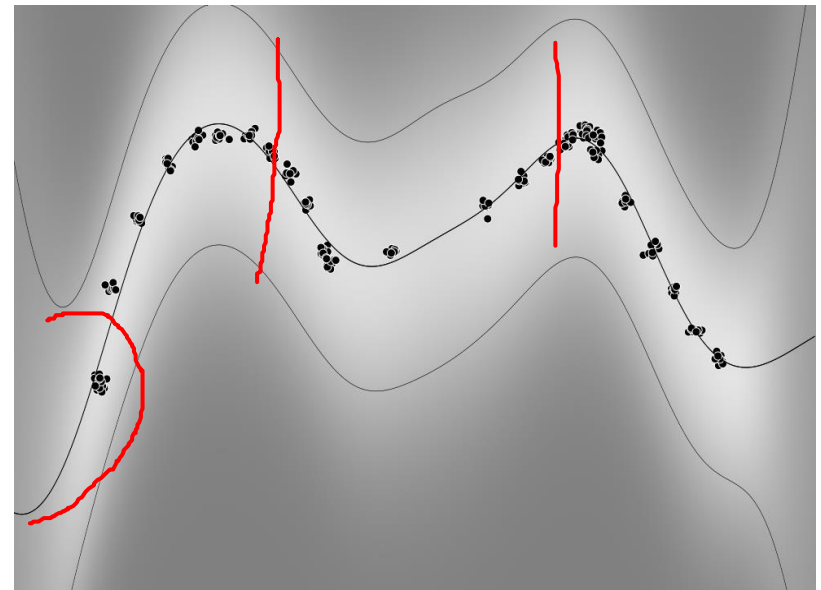
# Gaussian Process Regression (GPR)

## Hyperparameters

The larger the noise, the more uncertainty. The noise is  $\leq 1$ .



Low noise:  $\sigma=0.05$



High noise:  $\sigma=0.2$



# Gaussian Process Regression (GPR)

## Hyperparameters Tuning

Hyper-parameters: kernel's parameters and noise variance

One can automatically tune these hyperparameters by either:

- Crossvalidation
- Minimizing marginal Likelihood:

$$-\log(\mathbf{y} | \mathbf{X}, \sigma_\varepsilon, l) = 0.5 \underbrace{(\mathbf{y}^T (\mathbf{K} + \sigma_\varepsilon I)^{-1} \mathbf{y})}_{\text{Fit}} + \underbrace{\log |\mathbf{K} + \sigma_\varepsilon I|}_{\text{Complexity}} + \underbrace{\frac{M}{2} \log 2\pi}_{\text{Normalization}}$$

Automatically provides trade-off between fit and complexity



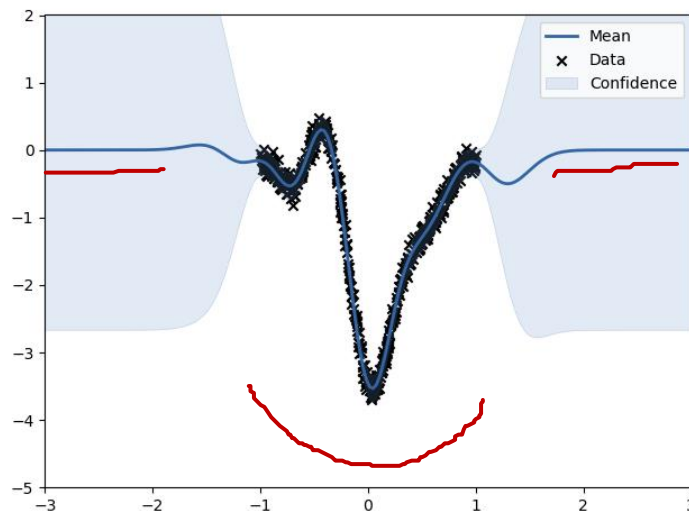
# Gaussian Process Regression (GPR)

## Prediction away from data

$$y = \sum_{i=1}^M \alpha_i \underline{k(x, x^i)}$$

$$\text{with } \alpha = \left[ K(X, X) + \sigma_\varepsilon^2 I \right]^{-1} \mathbf{y}$$

**GPR with RBF kernel predicts  $y=0$  away from datapoints!**



**Contrasts to SVR that predicts  $y=b$  away from datapoints.**



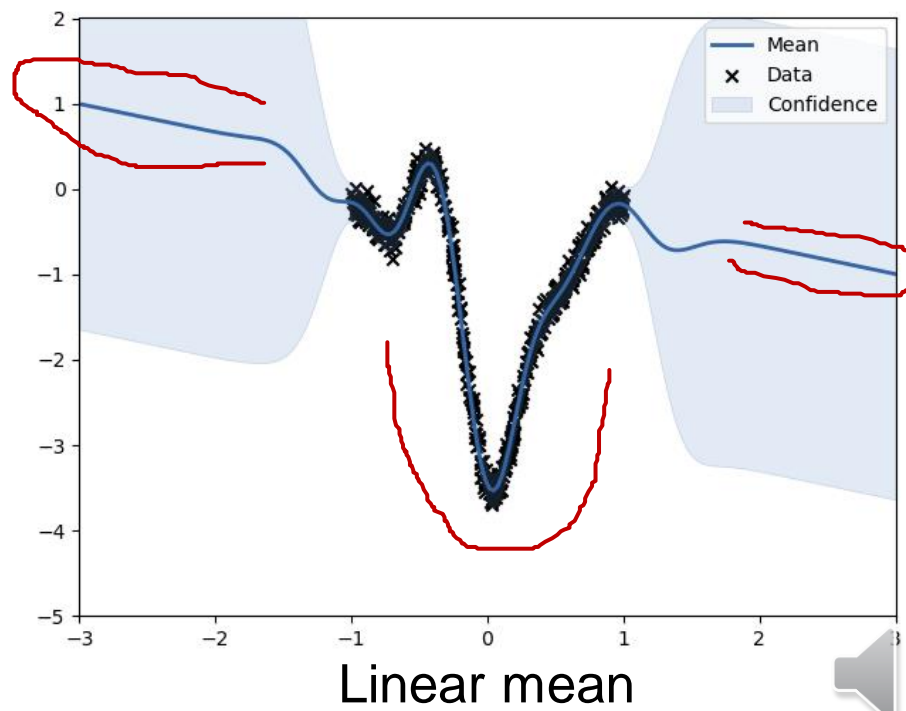
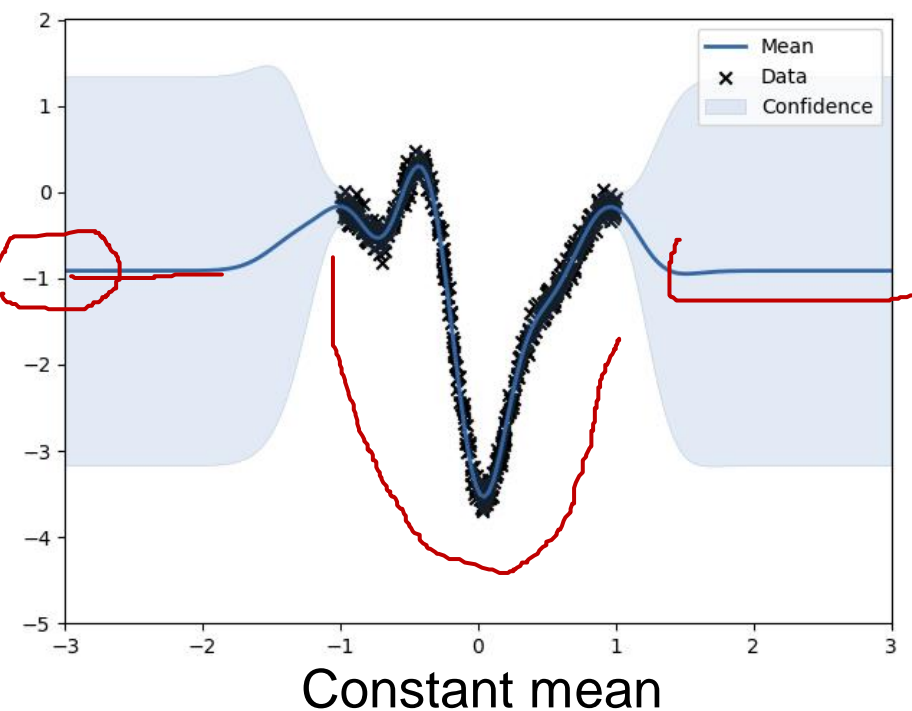
# Gaussian Process Regression (GPR)

## Prediction away from data

Instead of assuming zero mean one can add a mean function

$$\mathbb{E}\{p(y|x, X, \mathbf{y})\} = \underline{m(x)} + \sum_{i=1}^M \alpha_i k(\mathbf{x}^i, \mathbf{x})$$

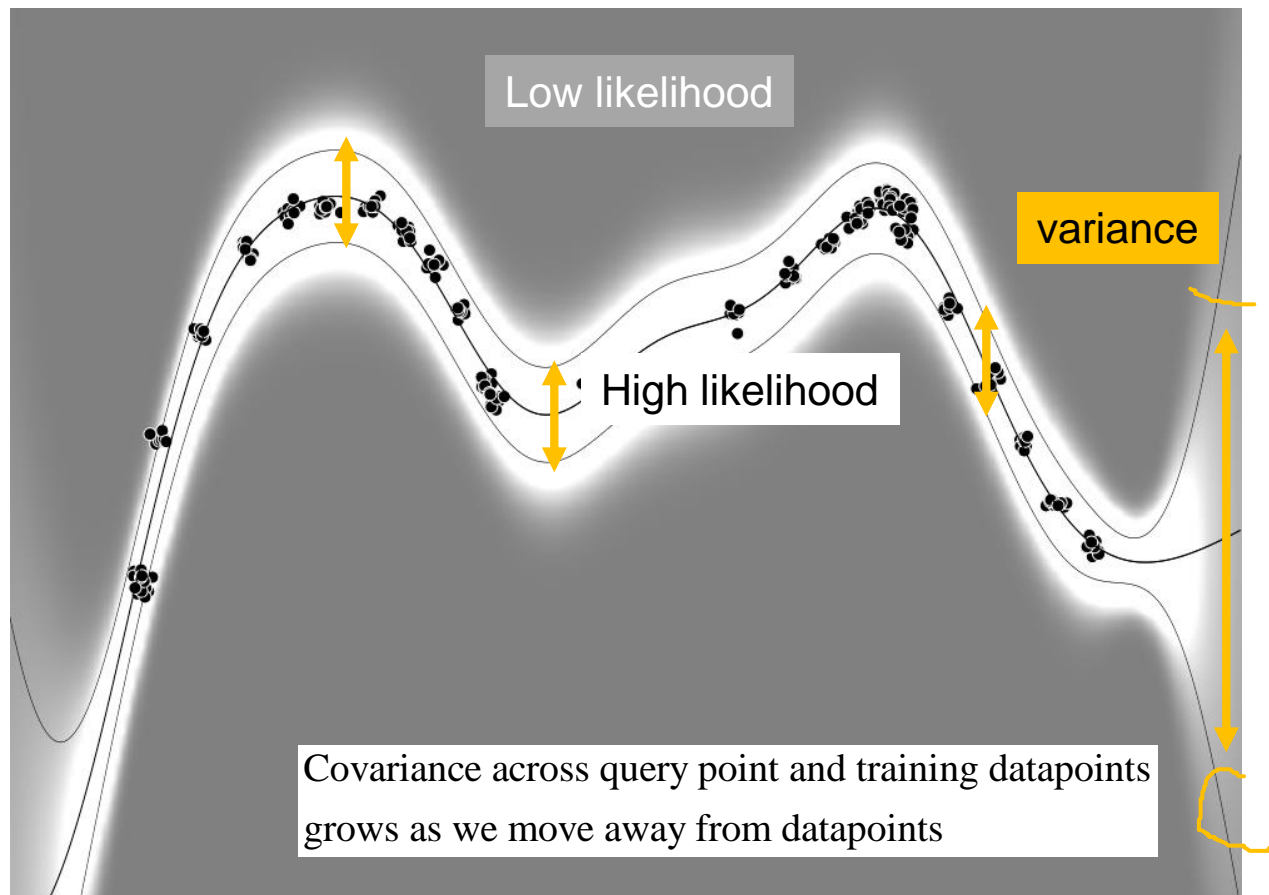
$$\text{where: } \underline{\alpha} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]^{-1} (\mathbf{y} - \underline{m(x)})$$



# Gaussian Process Regression (GPR)

## Confidence

The variance and the likelihood



# Gaussian Process Regression: Summary

## Advantages:

- ☐ Accuracy —
- ☐ Estimation of predictions' uncertainty —
- ☐ Auto-tuning of hyper-parameters —

## Disadvantage:

- ☐ Computational complexity  $O(M^3)$   
(but there exist several sparse methods for GP)

