

# CS-E4895 Gaussian Processes

## Lecture 2: Bayesian regression

Ti John

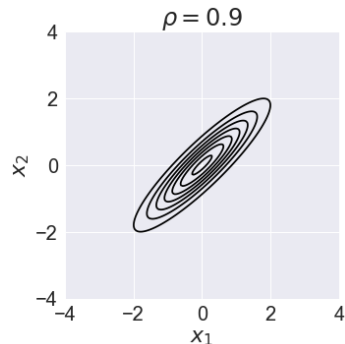
Aalto University

Tuesday 28.2.2023

# Last session

Last time, we talked about

- Course practicalities
- The multivariate Gaussian distribution
- The interpretation of the parameters
- Marginalization
- Conditional distributions
- How to sample from the distribution



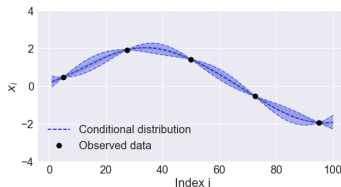
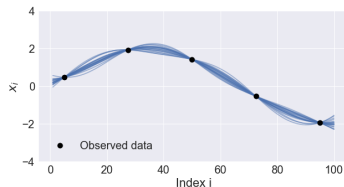
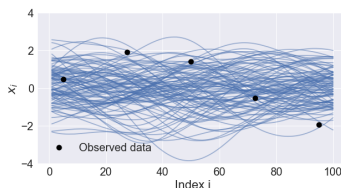
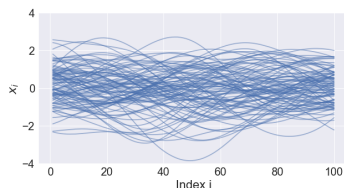
# Conditioning one more time

- Let  $\mathbf{x}_1$  and  $\mathbf{x}_2$  be a partitioning of  $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2$ , then

$$p(\mathbf{x}) = p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right) \quad (1)$$

- The conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2$  is:

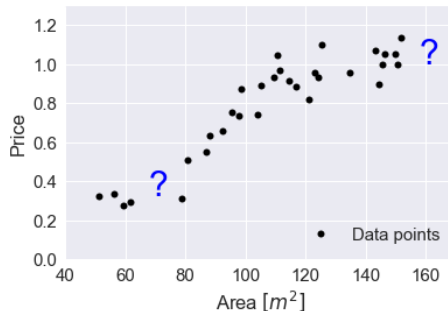
$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \Sigma_{12} \Sigma_{22}^{-1} [\mathbf{x}_2 - \mathbf{m}_2] + \mathbf{m}_1, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}) \quad (2)$$



# Gaussian processes for regression

## Running example

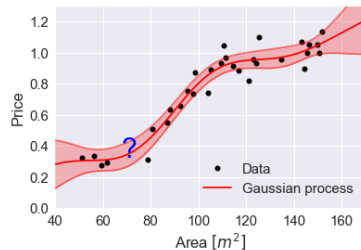
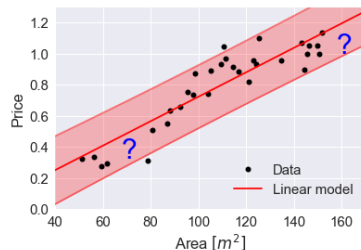
- Suppose we are given a data set of house prices in Helsinki



- Goal: Build a model using the data set and predict the average price for a house of  $70 m^2$  and  $160 m^2$

# Road map for today

- 1 The Bayesian linear model
- 2 The linear model as special case of a Gaussian process
- 3 Gaussian processes: definition & properties
- 4 Questions



# General setup for linear regression

- We are given a data set:  $\mathcal{D} = \{x_n, y_n\}_{n=1}^N$
- House example:  $y_n$  = house price and  $x_n$  = house area
- Goal: Learn some function  $f$  such that

$$y_n = f(\mathbf{x}_n) + \epsilon_n \quad (3)$$

- Assuming  $f$  is a linear model:

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_Dx_D = \sum_i w_ix_i = \mathbf{w}^\top \mathbf{x} \quad (4)$$

- Linear models are linear w.r.t. parameters, not the data:

$$f(\mathbf{x}) = w_1\phi_1(x_1) + w_2\phi_2(x_2) + \dots + w_{D'}\phi_{D'}(x_{D'}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}), \quad (5)$$

where  $\phi_i(\cdot)$  can be non-linear **feature** functions.

Which of the following models are linear models and why?

$$f(\mathbf{x}) = w_1 x_1 + w_2 x_2^2 + w_3 \sin(x_3) \quad (\text{Model 1})$$

$$f(\mathbf{x}) = w_1 x_1 + w_2^2 x_2 + w_3^3 x_3 \quad (\text{Model 2})$$

$$f(\mathbf{x}) = (\mathbf{w}^\top \mathbf{x})^2 \quad (\text{Model 3})$$

$$f(\mathbf{x}) = w_1 \exp(x_1) + w_2 \sqrt{x_2} + w_3 \quad (\text{Model 4})$$

Which of the following models are linear models and why?

$$f(\mathbf{x}) = w_1x_1 + w_2x_2^2 + w_3 \sin(x_3) \quad (\text{Model 1})$$

Yes,  $w_1$ ,  $w_2$ , and  $w_3$  all appear linearly

$$f(\mathbf{x}) = w_1x_1 + w_2^2x_2 + w_3^3x_3 \quad (\text{Model 2})$$

No, because of  $w_2^2$  and  $w_3^3$

$$f(\mathbf{x}) = (\mathbf{w}^\top \mathbf{x})^2 \quad (\text{Model 3})$$

No, because the weights will not appear linearly

$$f(\mathbf{x}) = w_1 \exp(x_1) + w_2 \sqrt{x_2} + w_3 \quad (\text{Model 4})$$

Yes,  $w_1$ ,  $w_2$ , and  $w_3$  all appear linearly



# Slope and intercept

- The models so far have not included an intercept or bias term
- Most often we want to incorporate an intercept/bias term

$$f(\mathbf{x}) = \textcolor{red}{w}_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \quad (6)$$

- By assuming  $x_0 = 1$ , we can write

$$\begin{aligned} f(\mathbf{x}) &= w_0 \cdot 1 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \\ &= w_0 \cdot x_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \\ &= \mathbf{w}^\top \mathbf{x} \end{aligned} \quad (7)$$

# Bayesian linear regression: Model and likelihood

- The model

$$y_n = f(\mathbf{x}_n) + \epsilon = \mathbf{w}^\top \mathbf{x}_n + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_{\text{obs}}^2) \quad (8)$$

- Likelihood for one data point

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(y_n | f(\mathbf{x}_n), \sigma_{\text{obs}}^2) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \sigma_{\text{obs}}^2) \quad (9)$$

- Likelihood for all data points

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{w}^\top \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma_{\text{obs}}^2 \mathbf{I}) \quad (10)$$

- Since the data is assumed constant, the likelihood is a function of parameters  $\mathbf{w}$
- Next step: we introduce a prior distribution  $p(\mathbf{w})$  for the weights  $\mathbf{w}$

# Bayesian linear regression: prior, posterior, evidence

- The prior  $p(\mathbf{w})$  contains our prior knowledge about  $\mathbf{w}$  **before** we see any data
- Bayes's rule gives us the posterior distribution

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \quad (11)$$

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w}) p(\mathbf{w})}{p(\mathbf{y})} \quad (12)$$

- Marginal likelihood (or *evidence*)

$$p(\mathbf{y}) = \int p(\mathbf{y}, \mathbf{w}) d\mathbf{w} = \int p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} = \mathbb{E}_{p(\mathbf{w})} [p(\mathbf{y}|\mathbf{w})]$$

- The posterior  $p(\mathbf{w}|\mathbf{y})$  captures everything we know about  $\mathbf{w}$  **after** seeing the data
- By convention we use  $p(\mathbf{w}|\mathbf{y})$  instead of the rigorous form  $p(\mathbf{w}|\mathbf{y}, \mathbf{X})$

# Bayesian linear regression: the posterior distribution

- We select a Gaussian prior for  $\mathbf{w}$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_p) \quad (13)$$

- The **parameter posterior** distribution becomes

$$p(\mathbf{w} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{w}) p(\mathbf{w})}{p(\mathbf{y})} \quad (14)$$

$$= \frac{\mathcal{N}(\mathbf{y} | \mathbf{X} \mathbf{w}, \sigma_{\text{obs}}^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_p)}{p(\mathbf{y})} \quad (15)$$

$$= \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \mathbf{A}^{-1}) \quad (16)$$

where

$$\boldsymbol{\mu} = \frac{1}{\sigma_{\text{obs}}^2} \mathbf{A}^{-1} \mathbf{X}^\top \mathbf{y} \quad \mathbf{A} = \frac{1}{\sigma_{\text{obs}}^2} \mathbf{X}^\top \mathbf{X} + \Sigma_p^{-1} \quad (17)$$

- See Rasmussen book section 2.1.1 for derivation (book eq 2.7).

# Bayesian linear regression: the predictive distribution

- We often want to compute the predictive distribution (or **predictive posterior**) for the noisy observation  $y_*$  at new data point  $\mathbf{x}_*$ , given as  $p(y_*|\mathbf{y})$
- We obtain the predictive distribution by averaging/marginalizing over the posterior:

$$p(y_*|\mathbf{y}) = \int p(y_*|\mathbf{x}_*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}) d\mathbf{w} \quad (18)$$

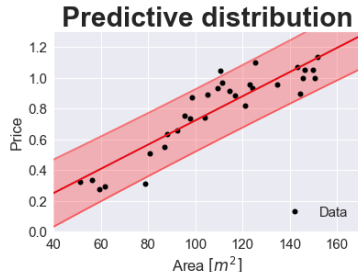
$$= \int \mathcal{N}(y_*|\mathbf{w}^\top \mathbf{x}_*, \sigma_{\text{obs}}^2) \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \mathbf{A}^{-1}) d\mathbf{w} \quad (19)$$

$$= \mathcal{N}(y_*|\boldsymbol{\mu}^\top \mathbf{x}_*, \sigma_{\text{obs}}^2 + \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*) \quad (20)$$

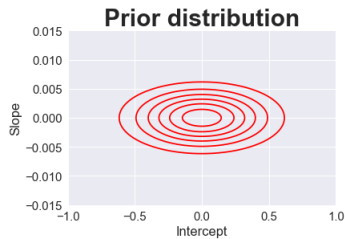
- The predictive distribution contains two sources of uncertainty:
  - ①  $\sigma_{\text{obs}}^2$ : measurement noise
  - ②  $\mathbf{A}^{-1}$ : uncertainty of the weights  $\mathbf{w}$
- $\mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*$ : uncertainty of the weights  $\mathbf{w}$  projected to the data space

# House price example: Posterior and predictive distributions

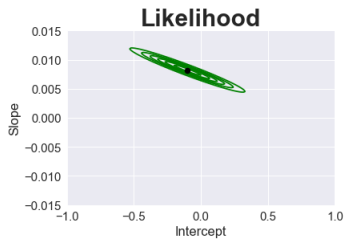
- The **posterior distribution** is a distribution over the parameter space
- The posterior is compromise between prior and likelihood
- The **predictive distribution** is a distribution over the output space



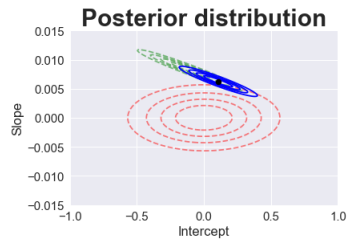
$$p(y^*|\mathbf{y}) = \mathcal{N}(y_* | \boldsymbol{\mu}^\top \mathbf{x}_*, \sigma_{\text{obs}}^2 + \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*)$$



$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \boldsymbol{\Sigma}_p)$$



$$p(\mathbf{y}|\mathbf{w}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma_{\text{obs}}^2 \mathbf{I})$$



$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \mathbf{A}^{-1})$$

Determine which of the following statements are true or false:

- ① Changing the prior distribution influences the posterior distribution
- ② Changing the prior distribution influences the likelihood
- ③ Changing the prior distribution influences the marginal likelihood
- ④ Changing the prior distribution influences the predictive distribution
- ⑤ The variance of the predictive distribution only depends on the measurement noise

Determine which of the following statements are true or false:

- ① Changing the prior distribution influences the posterior distribution  
true
- ② Changing the prior distribution influences the likelihood  
false
- ③ Changing the prior distribution influences the marginal likelihood  
true
- ④ Changing the prior distribution influences the predictive distribution  
true
- ⑤ The variance of the predictive distribution only depends on the measurement noise  
false



# Switching focus from parameters to functions (I)

- Our goal is to learn the function  $f$

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \quad (21)$$

- Until now we have focused on the weights  $\mathbf{w}$

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) \quad (22)$$

- Let's introduce  $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)] \in \mathbb{R}^N$  to the model  
The vector of predicted function values is  $\mathbf{f} = \mathbf{X}\mathbf{w}$

$$p(\mathbf{y}, \mathbf{f}, \mathbf{w}) = p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) \quad (23)$$

- Our model is still the same

$$p(\mathbf{y}, \mathbf{w}) = \int p(\mathbf{y}, \mathbf{f}, \mathbf{w}) \mathrm{d}\mathbf{f} = p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) \quad (24)$$

## Switching focus from parameters to functions (II)

- The augmented model

$$p(\mathbf{y}, \mathbf{f}, \mathbf{w}) = p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) \quad (25)$$

- What if we now marginalize over the weights?

$$p(\mathbf{y}, \mathbf{f}) = \int p(\mathbf{y}, \mathbf{f}, \mathbf{w}) d\mathbf{w} = p(\mathbf{y}|\mathbf{f}) \underbrace{\int p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w}}_{p(\mathbf{f})} \quad (26)$$

- We can decompose as likelihood and prior

$$p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}) \quad (27)$$

where

$$p(\mathbf{f}) = \int p(\mathbf{f}, \mathbf{w}) d\mathbf{w} = \int p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} \quad (28)$$

## Switching focus from parameters to functions (III)

- Let's study the prior distribution on  $\mathbf{f}$

$$p(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} = \int p(\mathbf{f}|\mathbf{w}) \mathcal{N}(\mathbf{w}|\mathbf{0}, \Sigma_p) d\mathbf{w} = ? \quad (29)$$

- We could do the integral directly...
- But let's instead use the result from last week

$$\mathbf{z} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \Rightarrow \quad \mathbf{A}\mathbf{z} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\mathbf{m} + \mathbf{b}, \mathbf{A}\mathbf{V}\mathbf{A}^\top) \quad (30)$$

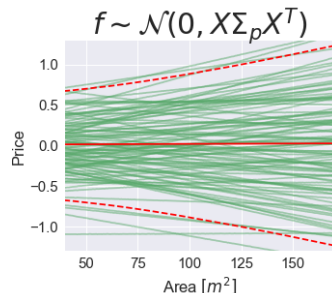
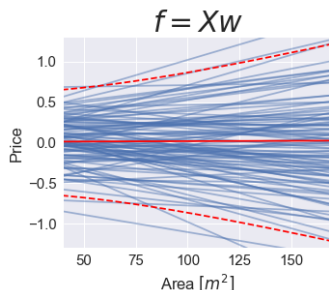
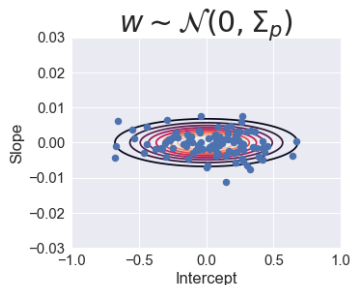
- We know that  $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \Sigma_p)$  and  $\mathbf{f} = \mathbf{X}\mathbf{w}$

$$\mathbb{E}[\mathbf{f}] = \mathbf{X}\mathbf{0} + \mathbf{0} = \mathbf{0} \qquad \mathbb{V}[\mathbf{f}] = \mathbf{X}\Sigma_p\mathbf{X}^\top \quad (31)$$

- In other words

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{X}\Sigma_p\mathbf{X}^\top) \quad (32)$$

# Weight view vs. function view



Same distribution for  $f$  in both cases but with two different representations

## Weight view

- Prior on weights:  $p(w)$
- $p(y, w) = p(y|w) p(w)$
- Posterior of weights:  $p(w|y)$

## Function view

- Prior on function values:  $p(f)$
- $p(y, f) = p(y|f) p(f)$
- Posterior of function values:  $p(f|y)$

## A closer look at the covariance matrix

- Prior on linear functions:  $p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$ , where  $\mathbf{K} = \mathbf{X}\Sigma_p\mathbf{X}^\top$
- Let's have a closer look at the covariance between  $f_i$  and  $f_j$

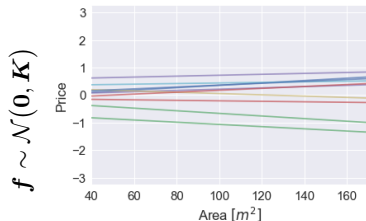
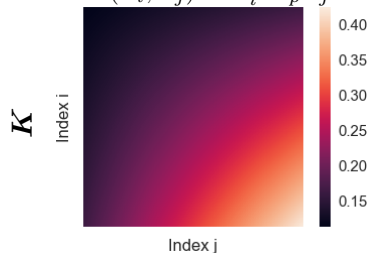
$$\begin{aligned}\mathbf{K}_{ij} &= \text{cov}(f_i, f_j) = \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j)) = \text{cov}(\mathbf{w}^\top \mathbf{x}_i, \mathbf{w}^\top \mathbf{x}_j) \\ &= \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_i - 0)(\mathbf{w}^\top \mathbf{x}_j - 0)] && \text{(Why zero mean?)} \\ &= \mathbb{E}[\mathbf{w}^\top \mathbf{x}_i \mathbf{w}^\top \mathbf{x}_j] \\ &= \mathbb{E}[\mathbf{x}_i^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_j] \\ &= \mathbf{x}_i^\top \mathbb{E}[\mathbf{w} \mathbf{w}^\top] \mathbf{x}_j \\ &= \mathbf{x}_i^\top \Sigma_p \mathbf{x}_j \\ &\equiv k(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

- The covariance function is called a **kernel** function
- What happens if we change the **covariance function**  $k(\mathbf{x}_i, \mathbf{x}_j)$ ?
- It would change  $f(\cdot)$  !

# Covariance functions

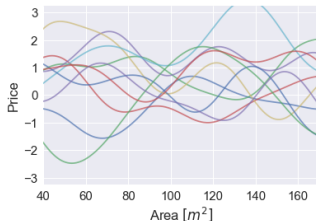
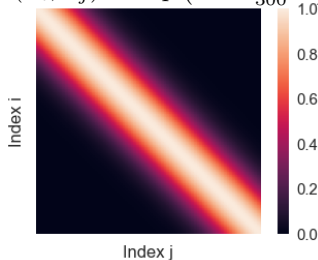
## Linear

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \Sigma_p \mathbf{x}_j$$



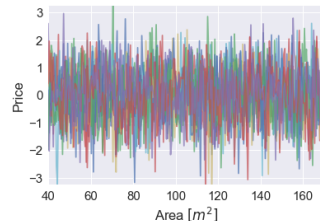
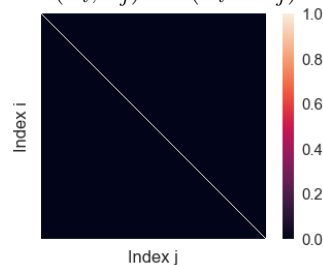
## Squared exponential

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|}{300}\right)$$



## White noise

$$k(\mathbf{x}_i, \mathbf{x}_j) = \delta(\mathbf{x}_i - \mathbf{x}_j)$$



The form of the covariance function determines the characteristics of functions

# Quiz

Consider the following covariance function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = 4 \quad \text{for all input pairs } (\mathbf{x}_i, \mathbf{x}_j) \quad (33)$$

- ① What is the marginal distribution of  $f(\mathbf{x}_i)$ ?
- ② What is the covariance between  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$ ?
- ③ What is the correlation between  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$ ?
- ④ What kind of functions are represented by the kernel in eq. (33)?

Consider the following covariance function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = 4 \quad \text{for all input pairs } (\mathbf{x}_i, \mathbf{x}_j) \quad (33)$$

- ① What is the marginal distribution of  $f(\mathbf{x}_i)$ ?  
 $p(f_i) = \mathcal{N}(0, 2^2) \quad (\sigma^2 = 4 \text{ or } \sigma = 2)$
- ② What is the covariance between  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$ ?  
 $\text{cov}(f_i, f_j) = 4$
- ③ What is the correlation between  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$ ?  
 $\text{corr}(f_i, f_j) = 1$
- ④ What kind of functions are represented by the kernel in eq. (33)?  
constant functions



# The big picture: Summary so far

- 1 We started with a Bayesian linear model

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w}) p(\mathbf{w}) \quad (34)$$

- 2 We introduced  $\mathbf{f}$  into the model and marginalized over the weights  $\mathbf{w}$

$$p(\mathbf{y}, \mathbf{f}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} = p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}) \quad (35)$$

- 3 This gave us a prior for linear functions in function space  $p(\mathbf{f})$ , where the covariance function for  $\mathbf{f}$  was given by

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \Sigma_p \mathbf{x} \quad (36)$$

- 4 By changing the form of the covariance function  $k(\mathbf{x}, \mathbf{x}')$ , we can model much more interesting functions

# Definitions

## Definition: multivariate Gaussian distribution

A random vector  $\mathbf{x} = [x_1, x_2, \dots, x_D]$  is said to have the **multivariate Gaussian distribution** if all linear combinations of  $\mathbf{x}$  are Gaussian distributed:

$$y = a_1x_1 + a_2x_2 + \dots + a_Dx_D \sim \mathcal{N}(m, v)$$

for all  $\mathbf{a} \in \mathbb{R}^D$

## Definition: Gaussian process

A **Gaussian process** is a collection of random variables indexed over space, any finite subset of which have a joint Gaussian distribution.

# Characterization and notation

- A Gaussian process can be considered as a prior distribution over functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  (the domain or index space  $\mathcal{X}$  is typically  $\mathbb{R}^D$ )

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (37)$$

- A Gaussian process is completely characterized by its mean function  $m(\mathbf{x})$  and its covariance function  $k(\mathbf{x}, \mathbf{x}')$ , which define

$$\mathbb{E}[f(\mathbf{x})] = m(\mathbf{x}) \quad (38)$$

$$\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \quad (39)$$

This means that  $f(\mathbf{x})$  and  $f(\mathbf{x}')$  are jointly Gaussian distributed with covariance  $k(\mathbf{x}, \mathbf{x}')$

- The probability of any subset of function **values**  $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$  at **any** inputs  $\mathbf{x}_1, \dots, \mathbf{x}_N$  is

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{m}, \mathbf{K}) \quad (40)$$

where  $\mathbf{m} = [m(\mathbf{x}_1), \dots, m(\mathbf{x}_N)]$  and  $[\mathbf{K}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$

# Gaussian processes are consistent wrt. marginalization

- Assume the function  $f$  follows a Gaussian process distribution:

$$f \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (41)$$

- The Gaussian process will induce a density for  $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2)]$ :

$$p(\mathbf{f}) = p(f_1, f_2) = \mathcal{N}\left(\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \mid \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix}\right) \quad (42)$$

- The induced density function for  $f_1 = f(\mathbf{x}_1)$  will always satisfy

$$p(f_1) = \mathcal{N}(f_1 \mid m_1, K_{11}) \quad (43)$$

- In words: “Examination of a larger set of variables does not change the distribution of the smaller set”
- If  $\mathcal{X} = \mathbb{R}^D$ , the GP prior describes infinitely many random variables  $\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D\}$ , but in practice we only have to deal with a finite subset corresponding to the data set at hand and where we want to evaluate or ‘test’ the function

# Gaussian process intuition

- A Gaussian process implements the assumption:

$$\mathbf{x} \approx \mathbf{x}' \Rightarrow f(\mathbf{x}) \approx f(\mathbf{x}') \quad (44)$$

- In other words: If the inputs are similar, the outputs should be similar as well.
- Using the squared exponential covariance function as example:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right) \quad (45)$$

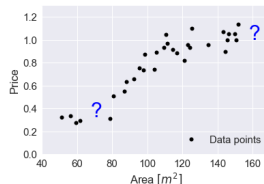
- Then covariance between  $f(\mathbf{x})$  and  $f(\mathbf{x}')$  is given by

$$\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right) \quad (46)$$

- Note: the covariance between **outputs** is given in terms of the **inputs**

## Back to our house price example (I)

**Goal:** To predict the price for a house with area  $x_* = 70 \text{ m}^2$  based on the training data  $\{x_n, y_n\}_{n=1}^N$



- Model:  $y_n = f(x_n)$ , where  $f$  is an unknown function (no noise for now)
- We impose a GP prior on  $f$ :  $\mathcal{GP}(m(x), k(x, x'))$ 
  - The prior is defined for all  $x \in \mathbb{R}$
  - We choose to evaluate the model at 70 observed points and evaluation points
- We choose  $m(x) = 0$  and the covariance function  $k(x, x')$  to be the squared exponential (and linear + bias term)
- The joint density for the training data becomes

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}_{ff}) \quad (47)$$

where  $\mathbf{f} = [f(x_1), f(x_2), \dots, f(x_N)]$  and  $[\mathbf{K}_{ff}]_{ij} = k(x_i, x_j)$

## Back to our house price example (II)

- The joint density for the training data

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}_{ff}) \quad (48)$$

- But what about the predictions for the new point  $x_*$  and the value of  $f(x_*)$ ?
- Let  $f_* = f(x_*)$ , then we can jointly model  $\mathbf{f}$  and  $f_*$  (consistency property)

$$p(\mathbf{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} | \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \quad (49)$$

where  $\mathbf{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_N)]^\top$  and  $K_{f_*f_*} = k(x_*, x_*)$

- Now we can use the rule for conditioning in Gaussian distributions to compute  $p(f_* | \mathbf{f})$

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top) \quad (50)$$

## Back to our house price example (III)

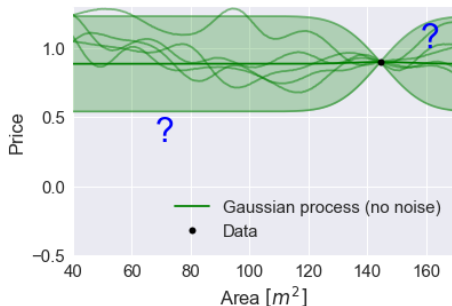
- The joint model for  $\mathbf{f}$  and  $f_*$  is

$$p(\mathbf{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

where  $\mathbf{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_N)]^\top$  and  $K_{f_*f_*} = k(x_*, x_*)$

- Conditioning on  $\mathbf{f}$  yields:

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$





## Back to our house price example (III)

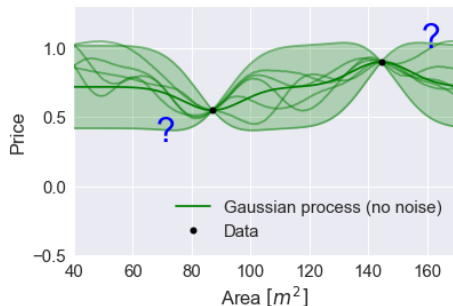
- The joint model for  $\mathbf{f}$  and  $f_*$  is

$$p(\mathbf{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

where  $\mathbf{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_N)]^\top$  and  $K_{f_*f_*} = k(x_*, x_*)$

- Conditioning on  $\mathbf{f}$  yields:

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$



## Back to our house price example (III)

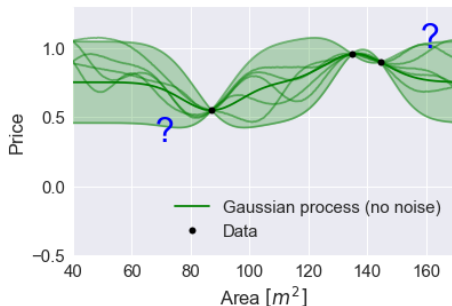
- The joint model for  $\mathbf{f}$  and  $f_*$  is

$$p(\mathbf{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

where  $\mathbf{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_N)]^\top$  and  $K_{f_*f_*} = k(x_*, x_*)$

- Conditioning on  $\mathbf{f}$  yields:

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$



## Back to our house price example (III)

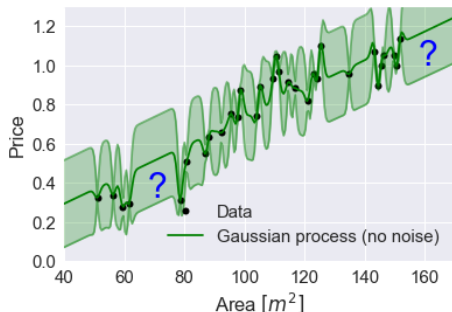
- The joint model for  $\mathbf{f}$  and  $f_*$  is

$$p(\mathbf{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

where  $\mathbf{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_N)]^\top$  and  $K_{f_*f_*} = k(x_*, x_*)$

- Conditioning on  $\mathbf{f}$  yields:

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$



## Back to our house price example (IV)

- Consider now the (more realistic) noisy model:  
 $y_n = f(x_n) + \epsilon_n$ , where  $\epsilon_n$  is Gaussian distributed
- Gaussian likelihood:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_{\text{obs}}^2 \mathbf{I}) \quad (51)$$

- The joint model for the noisy case becomes

$$\begin{aligned} p(\mathbf{y}, \mathbf{f}, f_*) &= p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) \\ &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_{\text{obs}}^2 \mathbf{I}) \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \end{aligned} \quad (52)$$

- Marginalizing over  $\mathbf{f}$  gives

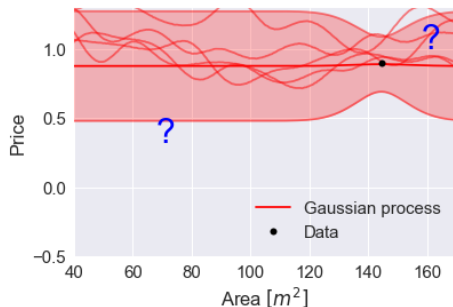
$$\begin{aligned} p(\mathbf{y}, f_*) &= \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f} \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \end{aligned} \quad (53)$$

## Back to our house price example (V)

- The joint distribution  $p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f}$   
$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top) \quad (54)$$

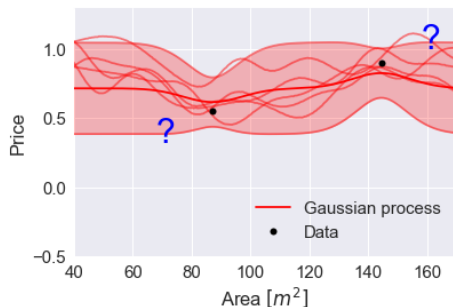


## Back to our house price example (V)

- The joint distribution  $p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f}$   
$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top) \quad (54)$$

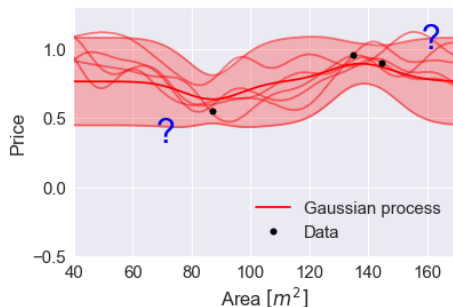


## Back to our house price example (V)

- The joint distribution  $p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f}$   
$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top) \quad (54)$$

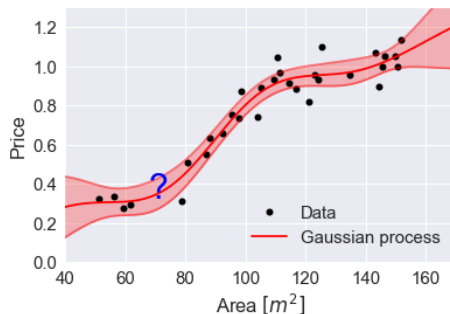


## Back to our house price example (V)

- The joint distribution 
$$p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f}$$
$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top) \quad (54)$$



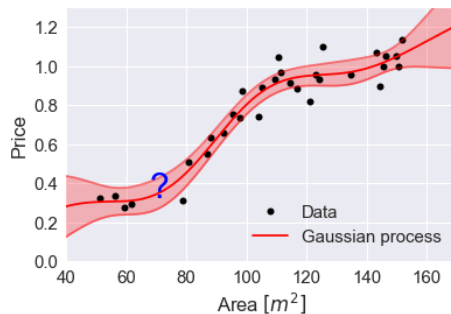
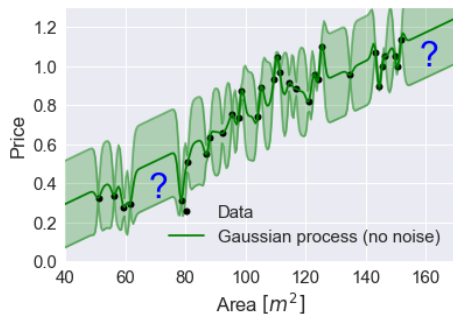


## Back to our house price example (V)

- The joint distribution  $p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}, f_*) d\mathbf{f}$ 
$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I} & \mathbf{K}_{ff_*} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right)$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top) \quad (54)$$



Posterior distribution in the noiseless case:

$$p(f_*|\mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$

Posterior distribution for the noisy case ( $y = f + \epsilon$ ):

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} (\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f} (\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top)$$

**Are the following statements true or false?:**

- ① Gaussian processes can fit highly non-linear functions, but the predictive means are given by a linear combination of the observations  $\mathbf{y}$ .
- ② The variance of the posterior distribution is independent of the observations  $\mathbf{y}$ .

Posterior distribution in the noiseless case:

$$p(f_*|\mathbf{f}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{f}, K_{f_*f_*} - \mathbf{K}_{f_*f} \mathbf{K}_{ff}^{-1} \mathbf{K}_{f_*f}^\top)$$

Posterior distribution for the noisy case ( $y = f + \epsilon$ ):

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f} (\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f} (\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top)$$

**Are the following statements true or false?:**

- 1 Gaussian processes can fit highly non-linear functions, but the predictive means are given by a linear combination of the observations  $\mathbf{y}$ .  
true
- 2 The variance of the posterior distribution is independent of the observations  $\mathbf{y}$ .  
true

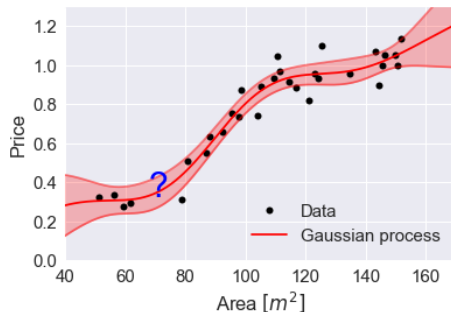
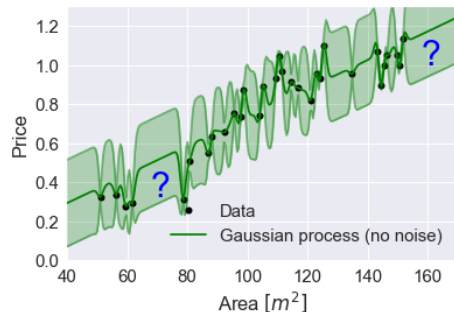
# What did we do?

- The predictive function posterior is conveniently a single equation (... for regression)

$$p(f_*|\mathbf{y}) = \mathcal{N}(f_* | \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K}_{f_*f_*} - \mathbf{K}_{f_*f}(\mathbf{K}_{ff} + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \mathbf{K}_{f_*f}^\top)$$

- We ended up not optimizing any parameters, how is this possible?
- Problem: how to define the hyperparameters
  - The noise variance  $\sigma_{\text{obs}}^2$
  - The kernel bandwidth or shape

⇒ Next lecture



# End of today's lecture

## Room change

- Exercise sessions: Thursdays 10:15 – 12:00, R001/U142 U4

## Next lecture:

- Kernels and covariance functions
- Model selection and hyperparameters
- Read ch. 4.2 and ch. 5.1-5.4 in Gaussian process book (<http://gaussianprocess.org/gpml/>)

## Assignment:

- Time to work on assignment #1  
(released on Wednesday, deadline next Wednesday 8th of March)
- Should be handed in (see Mycourses)
- In Jupyter notebook format