# CS-E4075 Special course on Gaussian processes: Session #2

Markus Heinonen

Aalto University
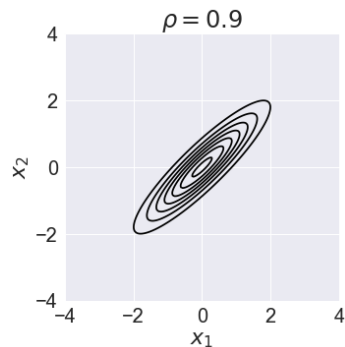
`markus.o.heinonen@aalto.fi`

Thursday 14.1.2021

# Last session

Last time, we talked about

- The multivariate Gaussian distribution

- The interpretation of the parameters

- Marginalization

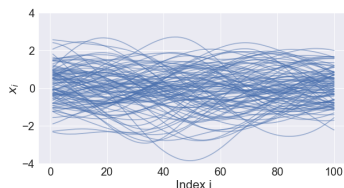- Conditional distributions

- How to sample from the distribution



$\rho = 0.9$

# Conditioning one more time

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

$$p(\boldsymbol{x}) = p(\boldsymbol{x}_1, \boldsymbol{x}_2) = \mathcal{N}\left( \begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \end{bmatrix} \Big| \begin{bmatrix} \boldsymbol{m}_1 \\ \boldsymbol{m}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right) \tag{1}$$

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

$$p(\boldsymbol{x}_1 | \boldsymbol{x}_2) = \mathcal{N}\left( \boldsymbol{x}_1 | \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \left[ \boldsymbol{x}_2 - \boldsymbol{m}_2 \right] + \boldsymbol{m}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \right) \tag{2}$$

# Conditioning one more time

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

$$p(\boldsymbol{x}) = p(\boldsymbol{x}_1, \boldsymbol{x}_2) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \end{bmatrix} \Big| \begin{bmatrix} \boldsymbol{m}_1 \\ \boldsymbol{m}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right) \qquad (1)$$

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

$$p(\boldsymbol{x}_1 | \boldsymbol{x}_2) = \mathcal{N}\left(\boldsymbol{x}_1 | \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} [\boldsymbol{x}_2 - \boldsymbol{m}_2] + \boldsymbol{m}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}\right) \qquad (2)$$

# Conditioning one more time

- Let $x_1$ and $x_2$ be a partitioning of $x = x_1 \cup x_2$, then

$$p(x) = p(x_1, x_2) = \mathcal{N}\left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \middle| \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right) \quad (1)$$

- The conditional distribution of $x_1$ is given $x_2$ by:

$$p(x_1|x_2) = \mathcal{N}\left( x_1 | \Sigma_{12}\Sigma_{22}^{-1}\left[ x_2 - m_2 \right] + m_1, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \right) \quad (2)$$

# Conditioning one more time

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

$$p(\boldsymbol{x}) = p(\boldsymbol{x}_1, \boldsymbol{x}_2) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \end{bmatrix} \Big| \begin{bmatrix} \boldsymbol{m}_1 \\ \boldsymbol{m}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right) \tag{1}$$

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

$$p(\boldsymbol{x}_1|\boldsymbol{x}_2) = \mathcal{N}\left(\boldsymbol{x}_1 | \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\left[\boldsymbol{x}_2 - \boldsymbol{m}_2\right] + \boldsymbol{m}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}\right) \tag{2}$$

# 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 $70m^2$ and $160m^2$

# 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

$$y_n = f(\boldsymbol{x}_n) + \epsilon_n \tag{3}$$

# 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

$$y_n = f(\boldsymbol{x}_n) + \epsilon_n \tag{3}$$

- Assuming $f$ is a linear model:

$$f(\boldsymbol{x}) = w_1 x_1 + w_2 x_2 + \ldots + w_D x_D = \sum_i w_i x_i = \boldsymbol{w}^T \boldsymbol{x} \tag{4}$$

# 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

$$y_n = f(\boldsymbol{x}_n) + \epsilon_n \tag{3}$$

- Assuming $f$ is a linear model:

$$f(\boldsymbol{x}) = w_1 x_1 + w_2 x_2 + \ldots + w_D x_D = \sum_i w_i x_i = \boldsymbol{w}^T \boldsymbol{x} \tag{4}$$

- Linear models are linear wrt. parameters, not the data:

$$f(\boldsymbol{x}) = w_1 \phi_1(x_1) + w_2 \phi_2(x_2) + \ldots + w_{D'} \phi_{D'}(x_{D'}) = \boldsymbol{w}^T \phi(\boldsymbol{x}), \tag{5}$$

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

## Question

Which of the following models are linear models and why?

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

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

$$f(\boldsymbol{x}) = \left(\boldsymbol{w}^T \boldsymbol{x}\right)^2 \tag{Model 3}$$

$$f(\boldsymbol{x}) = w_1 \exp(x_1) + w_2 \sqrt{x_2} + w_3 \tag{Model 4}$$

$$f(\boldsymbol{x}) = w_1 x_1 + w_2^2 x_2^2 + w_3^3 x_3^3 \tag{Model 5}$$

# 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

$$f(\boldsymbol{x}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots w_D x_D \tag{6}$$

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

$$f(\boldsymbol{x}) = w_0 \cdot 1 + w_1 x_1 + w_2 x_2 + \ldots w_D x_D \tag{7}$$

$$= w_0 \cdot x_0 + w_1 x_1 + w_2 x_2 + \ldots w_D x_D \tag{8}$$

$$= \boldsymbol{w}^T \boldsymbol{x} \tag{9}$$

# Bayesian linear regression

- The model

$$y_n = f(\boldsymbol{x}_n) + \epsilon = \boldsymbol{w}^T \boldsymbol{x}_n + \epsilon, \qquad \epsilon \sim \mathcal{N}\left(0, \sigma_{obs}^2\right) \tag{10}$$

# Bayesian linear regression

- The model

$$y_n = f(\mathbf{x}_n) + \epsilon = \mathbf{w}^T \mathbf{x}_n + \epsilon, \qquad \epsilon \sim \mathcal{N}\left(0, \sigma_{obs}^2\right) \tag{10}$$

- Likelihood for one data point

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(y_n | f(\mathbf{x}_n), \sigma_{obs}^2\right) = \mathcal{N}\left(y_n | \mathbf{w}^T \mathbf{x}_n, \sigma_{obs}^2\right) \tag{11}$$

# Bayesian linear regression

- The model

$$y_n = f(\mathbf{x}_n) + \epsilon = \mathbf{w}^T \mathbf{x}_n + \epsilon, \qquad \epsilon \sim \mathcal{N}\left(0, \sigma_{obs}^2\right) \tag{10}$$

- Likelihood for one data point

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(y_n | f(\mathbf{x}_n), \sigma_{obs}^2\right) = \mathcal{N}\left(y_n | \mathbf{w}^T \mathbf{x}_n, \sigma_{obs}^2\right) \tag{11}$$

- Likelihood for all data points

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^{N} p(y_n | \mathbf{w}^T \mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma_{obs}^2 \mathbf{I}\right) \tag{12}$$

# Bayesian linear regression

- The model

$$y_n = f(\boldsymbol{x}_n) + \epsilon = \boldsymbol{w}^T \boldsymbol{x}_n + \epsilon, \qquad \epsilon \sim \mathcal{N}\left(0, \sigma_{obs}^2\right) \tag{10}$$

- Likelihood for one data point

$$p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) = \mathcal{N}\left(y_n\big|f(\boldsymbol{x}_n), \sigma_{obs}^2\right) = \mathcal{N}\left(y_n\big|\boldsymbol{w}^T \boldsymbol{x}_n, \sigma_{obs}^2\right) \tag{11}$$

- Likelihood for all data points

$$p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w}) = \prod_{n=1}^{N} p(y_n|\boldsymbol{w}^T \boldsymbol{x}_n, \boldsymbol{w}) = \mathcal{N}\left(\boldsymbol{y}\big|\boldsymbol{X}\boldsymbol{w}, \sigma_{obs}^2 \boldsymbol{I}\right) \tag{12}$$

- Since the data is assumed constant, the likelihood is a function of parameters $\boldsymbol{w}$

- The prediction vector $\boldsymbol{f} = \boldsymbol{X}\boldsymbol{w}$

# Bayesian linear regression

- The model

$$y_n = f(\mathbf{x}_n) + \epsilon = \mathbf{w}^T \mathbf{x}_n + \epsilon, \qquad \epsilon \sim \mathcal{N}\left(0, \sigma_{obs}^2\right) \tag{10}$$

- Likelihood for one data point

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(y_n | f(\mathbf{x}_n), \sigma_{obs}^2\right) = \mathcal{N}\left(y_n | \mathbf{w}^T \mathbf{x}_n, \sigma_{obs}^2\right) \tag{11}$$

- Likelihood for all data points

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^{N} p(y_n | \mathbf{w}^T \mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma_{obs}^2 \mathbf{I}\right) \tag{12}$$

- Since the data is assumed constant, the likelihood is a function of parameters $\mathbf{w}$

- The prediction vector $\mathbf{f} = \mathbf{X}\mathbf{w}$

- Next step: we introduce a prior distribution $p(\mathbf{w})$ for the weights $\mathbf{w}$

# Bayesian linear regression

- The prior $p(\boldsymbol{w})$ contains our prior knowledge about $\boldsymbol{w}$ **before** we see any data

- Bayes rule gives us the posterior distribution

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \tag{13}$$

$$p(\boldsymbol{w}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y})} \tag{14}$$

# Bayesian linear regression

- The prior $p(\boldsymbol{w})$ contains our prior knowledge about **$\boldsymbol{w}$ before** we see any data

- Bayes rule gives us the posterior distribution

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \tag{13}$$

$$p(\boldsymbol{w}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y})} \tag{14}$$

- Marginal likelihood (or *evidence*)

$$p(\boldsymbol{y}) = \int p(\boldsymbol{y}, \boldsymbol{w}) \mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w}) \mathrm{d}\boldsymbol{w} = \mathbb{E}_{p(\boldsymbol{w})} p(\boldsymbol{y}|\boldsymbol{w})$$

# Bayesian linear regression

- The prior $p(\boldsymbol{w})$ contains our prior knowledge about $\boldsymbol{w}$ **before** we see any data

- Bayes rule gives us the posterior distribution

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \tag{13}$$

$$p(\boldsymbol{w}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y})} \tag{14}$$

- Marginal likelihood (or *evidence*)

$$p(\boldsymbol{y}) = \int p(\boldsymbol{y}, \boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \mathbb{E}_{p(\boldsymbol{w})}p(\boldsymbol{y}|\boldsymbol{w})$$

- The posterior $p(\boldsymbol{w}|\boldsymbol{y})$ captures everything we know about $\boldsymbol{w}$ **after** seing the data

- By convention we use $p(\boldsymbol{w}|\boldsymbol{y})$ instead of the rigorous form $p(\boldsymbol{w}|\boldsymbol{y}, \boldsymbol{X})$

# Bayesian linear regression: the posterior distribution

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

$$p(\boldsymbol{w}) = \mathcal{N}\left(\boldsymbol{w}\big|0, \boldsymbol{\Sigma}_p\right) \tag{15}$$

# Bayesian linear regression: the posterior distribution

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

$$p(\boldsymbol{w}) = \mathcal{N}\left(\boldsymbol{w}\big|0, \boldsymbol{\Sigma}_p\right) \tag{15}$$

- The parameter posterior distribution becomes

$$p(\boldsymbol{w}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y})} \tag{16}$$

$$= \frac{\mathcal{N}\left(\boldsymbol{y}\big|\boldsymbol{X}\boldsymbol{w}, \sigma_{obs}^2\boldsymbol{I}\right)\mathcal{N}\left(\boldsymbol{w}\big|0, \boldsymbol{\Sigma}_p\right)}{p(\boldsymbol{y})} \tag{17}$$

$$= \mathcal{N}\left(\boldsymbol{w}\big|\boldsymbol{\mu}, \boldsymbol{A}^{-1}\right) \tag{18}$$

where

$$\boldsymbol{\mu} = \frac{1}{\sigma_{obs}^2}\boldsymbol{A}^{-1}\boldsymbol{X}^T\boldsymbol{y} \qquad\qquad \boldsymbol{A} = \frac{1}{\sigma_{obs}^2}\boldsymbol{X}^T\boldsymbol{X} + \boldsymbol{\Sigma}_p^{-1} \tag{19}$$

- 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 $\boldsymbol{x}_*$, given as $p(y_*|\boldsymbol{y})$

# 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:

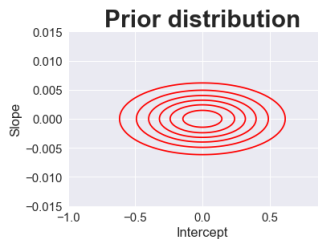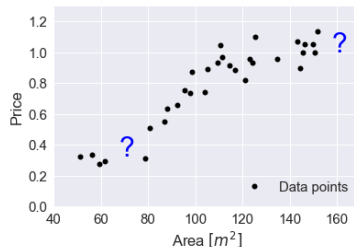$$p(y_*|\mathbf{y}) = \int p(y_*|\mathbf{x}_*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}) \mathrm{d}\mathbf{w} \tag{20}$$

$$= \int \mathcal{N}\left(y_*|\mathbf{w}^T\mathbf{x}_*, \sigma_{obs}^2\right) \mathcal{N}\left(\mathbf{w}|\boldsymbol{\mu}, \mathbf{A}^{-1}\right) \mathrm{d}\mathbf{w} \tag{21}$$

$$= \mathcal{N}\left(y_*|\boldsymbol{\mu}^T\mathbf{x}_*, \sigma_{obs}^2 + \mathbf{x}_*^T\mathbf{A}^{-1}\mathbf{x}_*\right) \tag{22}$$
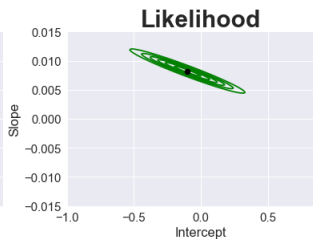
# Bayesian linear regression: the predictive distribution

- We often want to compute the predictive distribution (or predictive posterior) for the noisy observation $y_*$ at new data point $\boldsymbol{x}_*$, given as $p(y_*|\boldsymbol{y})$

- We obtain the predictive distribution by averaging/marginalizing over the posterior:

$$p(y_*|\boldsymbol{y}) = \int p(y_*|\boldsymbol{x}_*, \boldsymbol{w}) p(\boldsymbol{w}|\boldsymbol{y}) \mathrm{d}\boldsymbol{w} \tag{20}$$

$$= \int \mathcal{N}\left(y_*|\boldsymbol{w}^T\boldsymbol{x}_*, \sigma_{obs}^2\right) \mathcal{N}\left(\boldsymbol{w}|\boldsymbol{\mu}, \boldsymbol{A}^{-1}\right) \mathrm{d}\boldsymbol{w} \tag{21}$$

$$= \mathcal{N}\left(y_*|\boldsymbol{\mu}^T\boldsymbol{x}_*, \sigma_{obs}^2 + \boldsymbol{x}_*^T\boldsymbol{A}^{-1}\boldsymbol{x}_*\right) \tag{22}$$

- The predictive distributions contains two sources of uncertainty:

  1. $\sigma_{obs}^2$: measurement noise

  2. $\boldsymbol{A}^{-1}$: uncertainty of the weights $\boldsymbol{w}$

- $\boldsymbol{x}_*^T\boldsymbol{A}^{-1}\boldsymbol{x}_*$: uncertainty of the weights $\boldsymbol{w}$ projected to the data space

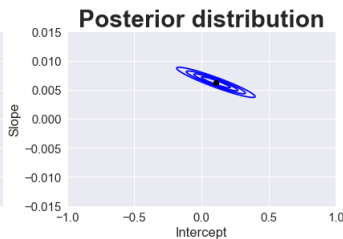# House price example: Posterior and predictive distributions

- The posterior distribution is distribution over the parameter space



**Prior distribution**

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

**Likelihood**

$$p(\boldsymbol{y}|\boldsymbol{w}) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{Xw}, \sigma_{obs}^2 \boldsymbol{I})$$

**Posterior distribution**

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

# House price example: Posterior and predictive distributions

- The posterior distribution is distribution over the parameter space

- The posterior is compromise between prior and likelihood





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

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

**Posterior distribution**
$$p(\boldsymbol{w}|\boldsymbol{y}) = \mathcal{N}(\boldsymbol{w}|\boldsymbol{\mu}, \boldsymbol{A}^{-1})$$
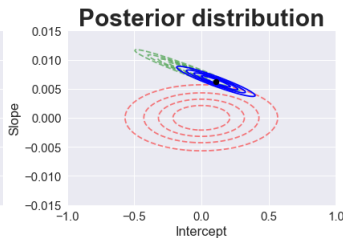
# House price example: Posterior and predictive distributions

- The posterior distribution is distribution over the parameter space

- The posterior is compromise between prior and likelihood

- The predictive distribution is a distribution over the output space



**Predictive distribution**

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



**Prior distribution**

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

**Likelihood**

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

**Posterior distribution**

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

# Question

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

- Our goal is to learn the function $f$

$$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} \tag{23}$$

# Switching focus from parameters to functions (I)

- Our goal is to learn the function $f$

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \tag{23}$$

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

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) \tag{24}$$

# Switching focus from parameters to functions (I)

- Our goal is to learn the function $f$

$$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} \tag{23}$$

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

$$p(\boldsymbol{y}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w}) \tag{24}$$

- Let's introduce $\boldsymbol{f} = [f(\boldsymbol{x}_1), f(\boldsymbol{x}_2), \ldots, f(\boldsymbol{x}_N)] \in \mathbb{R}^N$ to the model

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w}) \tag{25}$$

# Switching focus from parameters to functions (I)

- Our goal is to learn the function $f$

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \tag{23}$$

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

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) \tag{24}$$

- Let's introduce $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \ldots, f(\mathbf{x}_N)] \in \mathbb{R}^N$ to the model

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

- Our model is still the same

$$p(\mathbf{y}, \mathbf{w}) = \int p(\mathbf{y}, \mathbf{f}, \mathbf{w})\mathrm{d}\mathbf{f} = p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) \tag{26}$$

# Switching focus from parameters to functions (II)

- The augmented model

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w}) \qquad (27)$$

# Switching focus from parameters to functions (II)

- The augmented model

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w}) \tag{27}$$

- What if we now marginalize over the weights

$$p(\boldsymbol{y}, \boldsymbol{f}) = \int p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w})\mathrm{d}\boldsymbol{w} = p(\boldsymbol{y}|\boldsymbol{f})\underbrace{\int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w}}_{p(\boldsymbol{f})} \tag{28}$$

# Switching focus from parameters to functions (II)

- The augmented model

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w}) \tag{27}$$

- What if we now marginalize over the weights

$$p(\boldsymbol{y}, \boldsymbol{f}) = \int p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{w})\mathrm{d}\boldsymbol{w} = p(\boldsymbol{y}|\boldsymbol{f}) \underbrace{\int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w}}_{p(\boldsymbol{f})} \tag{28}$$

- We can decompose as likelihood and prior

$$p(\boldsymbol{y}, \boldsymbol{f}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}) \tag{29}$$

where

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}, \boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} \tag{30}$$

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} = ? \tag{31}$$

- We could do the integral directly...

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} = ? \tag{31}$$

- We could do the integral directly...

- But let's instead use the result from last week

$$\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{V}\right) \quad \Rightarrow \quad \boldsymbol{Az} + \boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{Am} + \boldsymbol{b}, \boldsymbol{AVA}^T\right) \tag{32}$$

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} = ? \tag{31}$$

- We could do the integral directly...

- But let's instead use the result from last week

$$\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{V}\right) \quad \Rightarrow \quad \boldsymbol{A}\boldsymbol{z} + \boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{A}\boldsymbol{m} + \boldsymbol{b}, \boldsymbol{A}\boldsymbol{V}\boldsymbol{A}^T\right) \tag{32}$$

- We know that $\boldsymbol{w} \sim \mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)$ and $\boldsymbol{f} = \boldsymbol{X}\boldsymbol{w}$

$$\mathbb{E}\left[\boldsymbol{f}\right] = \qquad\qquad\qquad \mathbb{V}\left[\boldsymbol{f}\right] = \tag{33}$$

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} =? \tag{31}$$

- We could do the integral directly...

- But let's instead use the result from last week

$$\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{V}\right) \quad \Rightarrow \quad \boldsymbol{A}\boldsymbol{z} + \boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{A}\boldsymbol{m} + \boldsymbol{b}, \boldsymbol{A}\boldsymbol{V}\boldsymbol{A}^T\right) \tag{32}$$

- We know that $\boldsymbol{w} \sim \mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)$ and $\boldsymbol{f} = \boldsymbol{X}\boldsymbol{w}$

$$\mathbb{E}\left[\boldsymbol{f}\right] = \boldsymbol{X}0 + 0 = 0 \qquad\qquad \mathbb{V}\left[\boldsymbol{f}\right] = \tag{33}$$

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} = ? \tag{31}$$

- We could do the integral directly...

- But let's instead use the result from last week

$$\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{V}\right) \quad \Rightarrow \quad \boldsymbol{Az} + \boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{Am} + \boldsymbol{b}, \boldsymbol{AVA}^T\right) \tag{32}$$

- We know that $\boldsymbol{w} \sim \mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)$ and $\boldsymbol{f} = \boldsymbol{Xw}$

$$\mathbb{E}\left[\boldsymbol{f}\right] = \boldsymbol{X}0 + 0 = 0 \qquad\qquad \mathbb{V}\left[\boldsymbol{f}\right] = \boldsymbol{X}\boldsymbol{\Sigma}_p\boldsymbol{X}^T \tag{33}$$

# Switching focus from parameters to functions (III)

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

$$p(\boldsymbol{f}) = \int p(\boldsymbol{f}|\boldsymbol{w})p(\boldsymbol{w})\mathrm{d}\boldsymbol{w} = \int p(\boldsymbol{f}|\boldsymbol{w})\mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)\mathrm{d}\boldsymbol{w} = ? \tag{31}$$

- We could do the integral directly...

- But let's instead use the result from last week

$$\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{V}\right) \quad \Rightarrow \quad \boldsymbol{Az} + \boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{Am} + \boldsymbol{b}, \boldsymbol{AVA}^T\right) \tag{32}$$

- We know that $\boldsymbol{w} \sim \mathcal{N}\left(\boldsymbol{w}|0, \boldsymbol{\Sigma}_p\right)$ and $\boldsymbol{f} = \boldsymbol{Xw}$

$$\mathbb{E}\left[\boldsymbol{f}\right] = \boldsymbol{X}0 + 0 = 0 \qquad\qquad \mathbb{V}\left[\boldsymbol{f}\right] = \boldsymbol{X}\boldsymbol{\Sigma}_p\boldsymbol{X}^T \tag{33}$$

- In other words

$$p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f}|0, \boldsymbol{X}\boldsymbol{\Sigma}_p\boldsymbol{X}^T\right) \tag{34}$$

# Weight view vs. function view

# Weight view vs. function view

# Weight view vs. function view

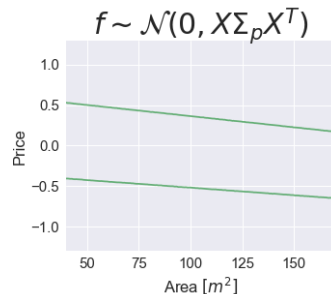# Weight view vs. function view
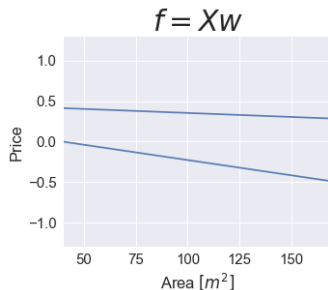
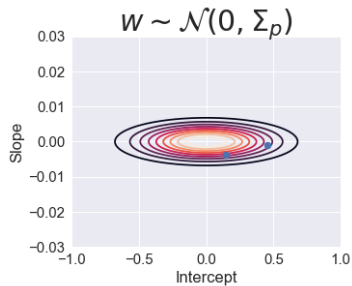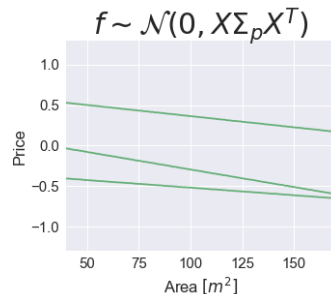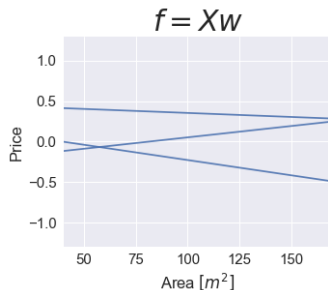# Weight view vs. function view

# Weight view vs. function view

# Weight view vs. function view

# Weight view vs. function view

# Weight view vs. function view

# Weight view vs. function view



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

**Weight view**

- Prior on weights: $p(\boldsymbol{w})$

- $p(\boldsymbol{y}, \boldsymbol{w}) = p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})$

- Posterior of weights: $p(\boldsymbol{w}|\boldsymbol{y})$

**Function view**

- Prior on function values: $p(\boldsymbol{f})$

- $p(\boldsymbol{y}, \boldsymbol{f}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f})$

- Posterior of function values: $p(\boldsymbol{f}|\boldsymbol{y})$

# A closer look at the covariance matrix

- Prior on linear functions: $p(\mathbf{f}) = \mathcal{N}\left(\mathbf{f}\,|\,0, \mathbf{K}\right)$, where $\mathbf{K} = \mathbf{X}\Sigma_p\mathbf{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

$$\mathbf{K}_{ij} = \mathrm{cov}\left(f_i, f_j\right) = \mathrm{cov}\left(f(\mathbf{x}_i), f(\mathbf{x}_j)\right) = \mathrm{cov}\left(\mathbf{w}^T\mathbf{x}_i, \mathbf{w}^T\mathbf{x}_j\right)$$

# A closer look at the covariance matrix

- Prior on linear functions: $p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f} | 0, \boldsymbol{K}\right)$, where $\boldsymbol{K} = \boldsymbol{X} \Sigma_p \boldsymbol{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

$$\begin{aligned}
\boldsymbol{K}_{ij} = \text{cov}\left(f_i, f_j\right) &= \text{cov}\left(f(\boldsymbol{x}_i), f(\boldsymbol{x}_j)\right) = \text{cov}\left(\boldsymbol{w}^T \boldsymbol{x}_i, \boldsymbol{w}^T \boldsymbol{x}_j\right) \\
&= \mathbb{E}\left[\left(\boldsymbol{w}^T \boldsymbol{x}_i - 0\right)\left(\boldsymbol{w}^T \boldsymbol{x}_j - 0\right)\right] \qquad \text{(Why zero mean?)}
\end{aligned}$$

# A closer look at the covariance matrix

- Prior on linear functions: $p(\mathbf{f}) = \mathcal{N}\left(\mathbf{f} | 0, \mathbf{K}\right)$, where $\mathbf{K} = \mathbf{X} \Sigma_p \mathbf{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

# A closer look at the covariance matrix

- Prior on linear functions: $p(\mathbf{f}) = \mathcal{N}\left(\mathbf{f}\,\middle|\,0, \mathbf{K}\right)$, where $\mathbf{K} = \mathbf{X}\Sigma_p\mathbf{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

# A closer look at the covariance matrix

- Prior on linear functions: $p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f}|0, \boldsymbol{K}\right)$, where $\boldsymbol{K} = \boldsymbol{X}\Sigma_p\boldsymbol{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

$$
\begin{aligned}
\boldsymbol{K}_{ij} = \text{cov}\left(f_i, f_j\right) &= \text{cov}\left(f(\boldsymbol{x}_i), f(\boldsymbol{x}_j)\right) = \text{cov}\left(\boldsymbol{w}^T\boldsymbol{x}_i, \boldsymbol{w}^T\boldsymbol{x}_j\right) \\
&= \mathbb{E}\left[\left(\boldsymbol{w}^T\boldsymbol{x}_i - 0\right)\left(\boldsymbol{w}^T\boldsymbol{x}_j - 0\right)\right] \qquad \text{(Why zero mean?)} \\
&= \mathbb{E}\left[\boldsymbol{w}^T\boldsymbol{x}_i\boldsymbol{w}^T\boldsymbol{x}_j\right] \\
&= \mathbb{E}\left[\boldsymbol{x}_i^T\boldsymbol{w}\boldsymbol{w}^T\boldsymbol{x}_j\right] \\
&= \boldsymbol{x}_i^T\mathbb{E}\left[\boldsymbol{w}\boldsymbol{w}^T\right]\boldsymbol{x}_j
\end{aligned}
$$

# A closer look at the covariance matrix

- Prior on linear functions: $p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f} | 0, \boldsymbol{K}\right)$, where $\boldsymbol{K} = \boldsymbol{X}\Sigma_p\boldsymbol{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

# A closer look at the covariance matrix

- Prior on linear functions: $p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f} \mid 0, \boldsymbol{K}\right)$, where $\boldsymbol{K} = \boldsymbol{X}\Sigma_p\boldsymbol{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

# A closer look at the covariance matrix

- Prior on linear functions: $p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f}|0, \boldsymbol{K}\right)$, where $\boldsymbol{K} = \boldsymbol{X}\Sigma_p\boldsymbol{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

- The covariance function is called a kernel function

- What happens if we change the **covariance function** $k\left(\boldsymbol{x}_i, \boldsymbol{x}_j\right)$?

# A closer look at the covariance matrix

- Prior on linear functions: $p(\mathbf{f}) = \mathcal{N}\left(\mathbf{f} \mid 0, \mathbf{K}\right)$, where $\mathbf{K} = \mathbf{X}\Sigma_p\mathbf{X}^T$

- Let's have a closer look on the covariance between $f_i$ and $f_j$

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

- The covariance function is called a kernel function

- What happens if we change the **covariance function** $k\left(\mathbf{x}_i, \mathbf{x}_j\right)$?

- It would change $f(\cdot)$ !

# Covariance functions



**Linear**
$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{\Sigma}_p \mathbf{x}_j$$

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

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

$\mathbf{K}$

Index i

Index i

Index i

Index j

Index j

Index j

$\mathbf{f} \sim \mathcal{N}(0, \mathbf{K})$

# Covariance functions

# Covariance functions



**Linear**

$$k\left(\boldsymbol{x}_i, \boldsymbol{x}_j\right) = \boldsymbol{x}_i^T \boldsymbol{\Sigma}_p \boldsymbol{x}_j$$

**Squared exponential**

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

**White noise**

$$k\left(\boldsymbol{x}_i, \boldsymbol{x}_j\right) = \delta\left(\boldsymbol{x}_i - \boldsymbol{x}_j\right)$$

# Covariance functions



**Linear**

$$k\left(\boldsymbol{x}_i, \boldsymbol{x}_j\right) = \boldsymbol{x}_i^T \boldsymbol{\Sigma}_p \boldsymbol{x}_j$$

**Squared exponential**

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

**White noise**

$$k\left(\boldsymbol{x}_i, \boldsymbol{x}_j\right) = \delta\left(\boldsymbol{x}_i - \boldsymbol{x}_j\right)$$

# Covariance functions



**Linear**
$$k(\boldsymbol{x}_i, \boldsymbol{x}_j) = \boldsymbol{x}_i^T \boldsymbol{\Sigma}_p \boldsymbol{x}_j$$

**Squared exponential**
$$k(\boldsymbol{x}_i, \boldsymbol{x}_j) = \exp\left(-\frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|}{300}\right)$$

**White noise**
$$k(\boldsymbol{x}_i, \boldsymbol{x}_j) = \delta(\boldsymbol{x}_i - \boldsymbol{x}_j)$$

The form of the covariance function determines the characteristics of functions

# Question

- Consider the following covariance function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = 1 \qquad \text{for all input pairs } (\mathbf{x}_i, \mathbf{x}_j) \tag{35}$$

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. (35)?

# The big picture: Summary so far

1. We started with a Bayesian linear model

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) \tag{36}$$

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

$$p(\mathbf{y}, \mathbf{f}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{w})p(\mathbf{w})\mathrm{d}\mathbf{w} = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}) \tag{37}$$

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

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{\Sigma}_p \mathbf{x} \tag{38}$$

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

$$y = a_1 x_1 + a_2 x_2 + \cdots + a_D x_D \sim \mathcal{N}(m, v)$$

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

## Definition: Gaussian process

A **Gaussian process** is a collection of random variables index 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 $f : \mathcal{X} \to \mathbb{R}$ (the domain or index space $\mathcal{X}$ is typically $\mathbb{R}^D$)

$$f(\boldsymbol{x}) \sim \mathcal{GP}\left(m\left(\boldsymbol{x}\right), k\left(\boldsymbol{x}, \boldsymbol{x}'\right)\right) \qquad (39)$$

# Characterization and notation

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

$$f(\boldsymbol{x}) \sim \mathcal{GP}\left(m\left(\boldsymbol{x}\right), k\left(\boldsymbol{x}, \boldsymbol{x}'\right)\right) \tag{39}$$

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

$$\mathbb{E}\left[f(\boldsymbol{x})\right] = m\left(\boldsymbol{x}\right) \tag{40}$$
$$cov[f(\boldsymbol{x}), f(\boldsymbol{x}')] = k\left(\boldsymbol{x}, \boldsymbol{x}'\right) \tag{41}$$

# Characterization and notation

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

$$f(\boldsymbol{x}) \sim \mathcal{GP}\left(m\left(\boldsymbol{x}\right), k\left(\boldsymbol{x}, \boldsymbol{x}'\right)\right) \tag{39}$$

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

$$\mathbb{E}\left[f(\boldsymbol{x})\right] = m\left(\boldsymbol{x}\right) \tag{40}$$
$$cov[f(\boldsymbol{x}), f(\boldsymbol{x}')] = k\left(\boldsymbol{x}, \boldsymbol{x}'\right) \tag{41}$$

- The probability of any subset of function values $\boldsymbol{f} = f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_N)$ at any inputs $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$ is

$$p(\boldsymbol{f}) = \mathcal{N}(\boldsymbol{f}|\boldsymbol{m}, \boldsymbol{K}) \tag{42}$$

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

# Gaussian processes are consistent wrt. marginalization

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

$$f \sim \mathcal{GP}\left(m\left(\mathbf{x}\right), k\left(\mathbf{x}, \mathbf{x}'\right)\right) \tag{43}$$

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

$$p(\mathbf{f}) = p(f_1, f_2) = \mathcal{N}\left(\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \Big| \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix}\right) \tag{44}$$

# Gaussian processes are consistent wrt. marginalization

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

$$f \sim \mathcal{GP}\left(m\left(\boldsymbol{x}\right), k\left(\boldsymbol{x}, \boldsymbol{x}'\right)\right) \tag{43}$$

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

$$p(\boldsymbol{f}) = p(f_1, f_2) = \mathcal{N}\left(\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \Bigg| \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix}\right) \tag{44}$$

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

$$p(f_1) = \mathcal{N}\left(f_1 | m_1, K_{11}\right) \tag{45}$$

- In words: "Examination of a larger set of variables does not change the distribution of the smaller set"

# Gaussian processes are consistent wrt. marginalization

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

$$f \sim \mathcal{GP}\left(m\left(\boldsymbol{x}\right), k\left(\boldsymbol{x}, \boldsymbol{x}'\right)\right) \tag{43}$$

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

$$p(\boldsymbol{f}) = p(f_1, f_2) = \mathcal{N}\left(\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \Big| \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix}\right) \tag{44}$$

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

$$p(f_1) = \mathcal{N}\left(f_1 | m_1, K_{11}\right) \tag{45}$$

- 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 variable $\left\{f(\boldsymbol{x}) : \boldsymbol{x} \in \mathbb{R}^D\right\}$, 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

- Gaussian process implements the assumption:

$$\mathbf{x} \approx \mathbf{x}' \quad \Rightarrow \quad f(\mathbf{x}) \approx f(\mathbf{x}') \tag{46}$$

- In other words: If the inputs are similar, the outputs should be similar as well.

# Gaussian process intuition

- Gaussian process implements the assumption:

$$\boldsymbol{x} \approx \boldsymbol{x}' \quad \Rightarrow \quad f(\boldsymbol{x}) \approx f(\boldsymbol{x}') \tag{46}$$

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

$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2}\right) \tag{47}$$

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

$$\mathrm{cov}\left[f(\boldsymbol{x}), f(\boldsymbol{x}')\right] = k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2}\right) \tag{48}$$

# Gaussian process intuition

- Gaussian process implements the assumption:

$$\boldsymbol{x} \approx \boldsymbol{x}' \quad \Rightarrow \quad f(\boldsymbol{x}) \approx f(\boldsymbol{x}') \tag{46}$$

- In other words: If the inputs are similar, the outputs should be similar as well.

- Using the squared exponential covariance function as example

$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2}\right) \tag{47}$$

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

$$\mathrm{cov}\left[f(\boldsymbol{x}), f(\boldsymbol{x}')\right] = k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2}\right) \tag{48}$$

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

# Back to our house price example (I)

**Goal**: To predict to the price for a house with area $x_* = 70$ 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 $f$: $\mathcal{GP}\left(m(x), k(x, x')\right)$
    - 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 $k(x, x')$ to be the covariance function to be the squared exponential (and linear + bias term)

- The joint density for the training data becomes

$$p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f} \middle| 0, \boldsymbol{K}_{ff}\right) \tag{49}$$

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

# Back to our house price example (II)

- The joint density for the training data

$$p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f}\,\middle|\,0, \boldsymbol{K}_{ff}\right) \tag{50}$$

- But what about the predictions for the new point $x_*$ and the value of $f(x_*)$?

# Back to our house price example (II)

- The joint density for the training data

$$p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f}\,\middle|\,0, \boldsymbol{K}_{ff}\right) \tag{50}$$

- 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 $\boldsymbol{f}$ and $f_*$ (consistency property)

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \middle| 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*} \\ \boldsymbol{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{51}$$

where $\boldsymbol{K}_{f_*f} = [k(x_*, x_1), k(x_*, x_2), \ldots, k(x_*, x_N)]$ and $K_{f_*f_*} = k(x_*, x_*)$

# Back to our house price example (II)

- The joint density for the training data

$$p(\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{f} \middle| 0, \boldsymbol{K}_{ff}\right) \tag{50}$$

- 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 $\boldsymbol{f}$ and $f_*$ (consistency property)

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix}\boldsymbol{f}\\f_*\end{bmatrix} \middle| 0, \begin{bmatrix}\boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*}\\\boldsymbol{K}_{f_*f} & K_{f_*f_*}\end{bmatrix}\right) \tag{51}$$

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

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

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_* \middle| \boldsymbol{K}_{f_*f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{K}_{f_*f}^T\right) \tag{52}$$

# Back to our house price example (III)

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

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \bigg| 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*} \\ \boldsymbol{K}_{f_* f} & K_{f_* f_*} \end{bmatrix}\right) \tag{53}$$

  where $\boldsymbol{K}_{f_* f} = [k(x_*, x_1), k(x_*, x_2), \ldots, k(x_*, x_N)]$ and $K_{f_* f_*} = k(x_*, x_*)$

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

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_* \big| \boldsymbol{K}_{f_* f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{y}, K_{f_* f_*} - \boldsymbol{K}_{f_* f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{K}_{f_* f}^T\right) \tag{54}$$

# Back to our house price example (III)

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

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \middle| 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*} \\ \boldsymbol{K}_{f_* f} & K_{f_* f_*} \end{bmatrix}\right) \tag{53}$$

  where $\boldsymbol{K}_{f_* f} = [k(x_*, x_1), k(x_*, x_2), \ldots, k(x_*, x_N)]$ and $K_{f_* f_*} = k(x_*, x_*)$

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

$$p(f_* | \boldsymbol{f}) = \mathcal{N}\left(f_* \middle| \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{y}, K_{f_* f_*} - \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{K}_{f_* f}^T\right) \tag{54}$$

# Back to our house price example (III)

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

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \Big| 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*} \\ \boldsymbol{K}_{f_* f} & K_{f_* f_*} \end{bmatrix}\right) \tag{53}$$

  where $\boldsymbol{K}_{f_* f} = [k(x_*, x_1), k(x_*, x_2), \ldots, k(x_*, x_N)]$ and $K_{f_* f_*} = k(x_*, x_*)$

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

$$p(f_* | \boldsymbol{f}) = \mathcal{N}\left(f_* \big| \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{y}, K_{f_* f_*} - \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{K}_{f_* f}^T\right) \tag{54}$$

# Back to our house price example (III)

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

$$p(\boldsymbol{f}, f_*) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \middle| 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{ff_*} \\ \boldsymbol{K}_{f_* f} & K_{f_* f_*} \end{bmatrix}\right) \tag{53}$$

  where $\boldsymbol{K}_{f_* f} = [k(x_*, x_1), k(x_*, x_2), \ldots, k(x_*, x_N)]$ and $K_{f_* f_*} = k(x_*, x_*)$

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

$$p(f_* | \boldsymbol{f}) = \mathcal{N}\left(f_* \middle| \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{y}, K_{f_* f_*} - \boldsymbol{K}_{f_* f} \boldsymbol{K}_{ff}^{-1} \boldsymbol{K}_{f_* f}^T\right) \tag{54}$$

# Back to our house price example (IV)

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

- Gaussian likelihood:

$$p(\boldsymbol{y}|\boldsymbol{f}) = \mathcal{N}\left(\boldsymbol{y}|\boldsymbol{f}, \sigma_{obs}^2 \boldsymbol{I}\right) \tag{55}$$

- The joint model for the noisy case becomes

$$p(\boldsymbol{y}, \boldsymbol{f}, f_*) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}, f_*) \tag{56}$$

$$= \mathcal{N}\left(\boldsymbol{y}|\boldsymbol{f}, \sigma_{obs}^2 \boldsymbol{I}\right) \mathcal{N}\left(\begin{bmatrix} \boldsymbol{f} \\ f_* \end{bmatrix} \boldsymbol{f} | 0, \begin{bmatrix} \boldsymbol{K}_{ff} & \boldsymbol{K}_{f_*f} \\ \boldsymbol{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{57}$$

# Back to our house price example (IV)

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

- Gaussian likelihood:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}\left(\mathbf{y}|\mathbf{f}, \sigma_{obs}^2 \mathbf{I}\right) \tag{55}$$

- The joint model for the noisy case becomes

$$p(\mathbf{y}, \mathbf{f}, f_*) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, f_*) \tag{56}$$

$$= \mathcal{N}\left(\mathbf{y}|\mathbf{f}, \sigma_{obs}^2 \mathbf{I}\right) \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \mathbf{f} | 0, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{f_*f} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{57}$$

- Marginalizing over $\mathbf{f}$ gives

$$p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, f_*)\mathrm{d}\mathbf{f} \tag{58}$$

$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \mathbf{f} | 0, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{obs}^2 \mathbf{I} & \mathbf{K}_{f_*f} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{59}$$
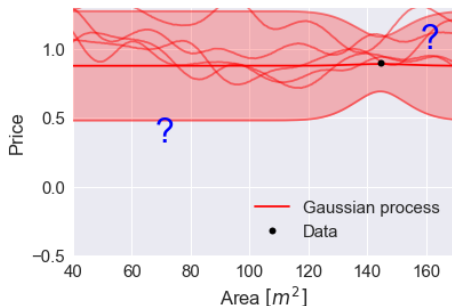
# Back to our house price example (V)

- The joint distribution

$$p(\boldsymbol{y}, f_*) = \int p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}, f_*)\mathrm{d}\boldsymbol{f} \tag{60}$$

$$= \mathcal{N}\left(\begin{bmatrix} \boldsymbol{y} \\ f_* \end{bmatrix}\Big|0, \begin{bmatrix} \boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I} & \boldsymbol{K}_{f_*f} \\ \boldsymbol{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{61}$$

- Once again, we can use the rule for conditioning

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_*\Big|\boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{K}_{f_*f}^T\right) \tag{62}$$
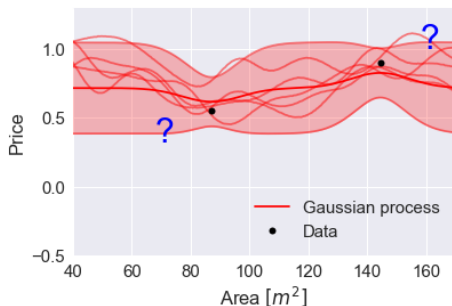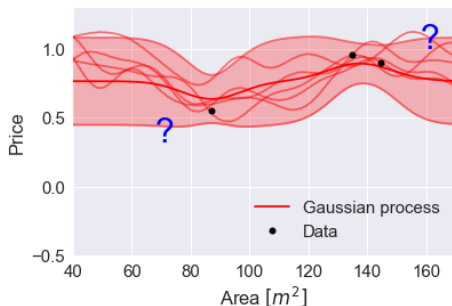
# Back to our house price example (V)

- The joint distribution

$$p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, f_*)\mathrm{d}\mathbf{f} \tag{60}$$

$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \Big| 0, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{obs}^2\mathbf{I} & \mathbf{K}_{f_*f} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \tag{61}$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{f}) = \mathcal{N}\left(f_*\big|\mathbf{K}_{f_*f}\left(\mathbf{K}_{ff} + \sigma_{obs}^2\mathbf{I}\right)^{-1}\mathbf{y}, K_{f_*f_*} - \mathbf{K}_{f_*f}\left(\mathbf{K}_{ff} + \sigma_{obs}^2\mathbf{I}\right)^{-1}\mathbf{K}_{f_*f}^T\right) \tag{62}$$
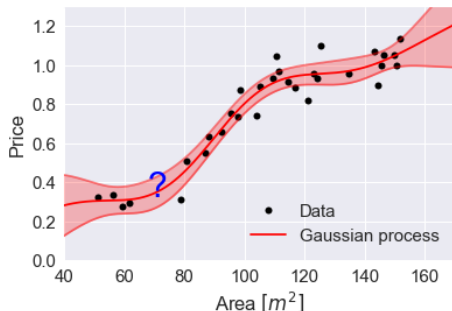
# Back to our house price example (V)

- The joint distribution

$$p(\mathbf{y}, f_*) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, f_*)\mathrm{d}\mathbf{f} \tag{60}$$

$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| 0, \begin{bmatrix} \mathbf{K}_{ff} + \sigma_{obs}^2 \mathbf{I} & \mathbf{K}_{f_* f} \\ \mathbf{K}_{f_* f} & K_{f_* f_*} \end{bmatrix}\right) \tag{61}$$

- Once again, we can use the rule for conditioning

$$p(f_*|\mathbf{f}) = \mathcal{N}\left(f_* \middle| \mathbf{K}_{f_* f}\left(\mathbf{K}_{ff} + \sigma_{obs}^2\mathbf{I}\right)^{-1}\mathbf{y}, K_{f_* f_*} - \mathbf{K}_{f_* f}\left(\mathbf{K}_{ff} + \sigma_{obs}^2\mathbf{I}\right)^{-1}\mathbf{K}_{f_* f}^T\right) \tag{62}$$
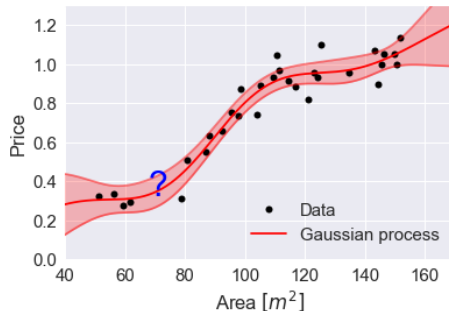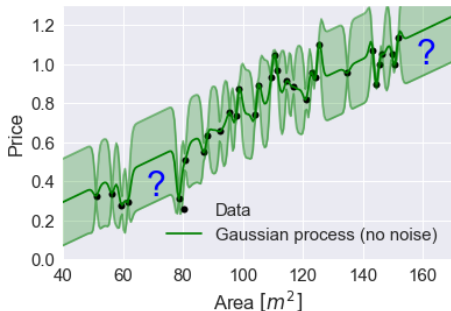
# Back to our house price example (V)

- The joint distribution

$$p(\boldsymbol{y}, f_*) = \int p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}, f_*)\mathrm{d}\boldsymbol{f} \tag{60}$$

$$= \mathcal{N}\left(\begin{bmatrix}\boldsymbol{y}\\f_*\end{bmatrix}\Big|0, \begin{bmatrix}\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I} & \boldsymbol{K}_{f_*f}\\\boldsymbol{K}_{f_*f} & K_{f_*f_*}\end{bmatrix}\right) \tag{61}$$

- Once again, we can use the rule for conditioning

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_*\big|\boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{K}_{f_*f}^T\right) \tag{62}$$

# Back to our house price example (V)

- The joint distribution

$$p(\boldsymbol{y}, f_*) = \int p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}, f_*)\mathrm{d}\boldsymbol{f} \tag{60}$$

$$= \mathcal{N}\left(\begin{bmatrix}\boldsymbol{y} \\ f_*\end{bmatrix}|0, \begin{bmatrix}\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I} & \boldsymbol{K}_{f_* f} \\ \boldsymbol{K}_{f_* f} & K_{f_* f_*}\end{bmatrix}\right) \tag{61}$$

- Once again, we can use the rule for conditioning

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_*|\boldsymbol{K}_{f_* f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{y}, K_{f_* f_*} - \boldsymbol{K}_{f_* f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{K}_{f_* f}^T\right) \tag{62}$$

## Question

Posterior distribution in the noiseless case:

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_*\big|\boldsymbol{K}_{f_*f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\boldsymbol{K}_{ff}^{-1}\boldsymbol{K}_{f_*f}^{T}\right) \tag{63}$$

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

$$p(f_*|\boldsymbol{y}) = \mathcal{N}\left(f_*\big|\boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{K}_{f_*f}^{T}\right) \tag{64}$$
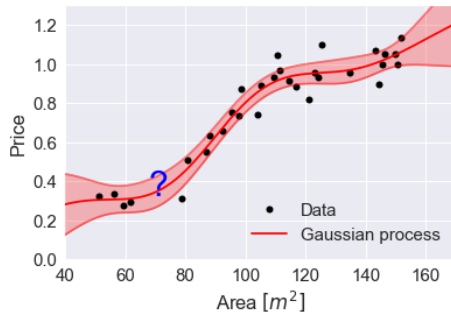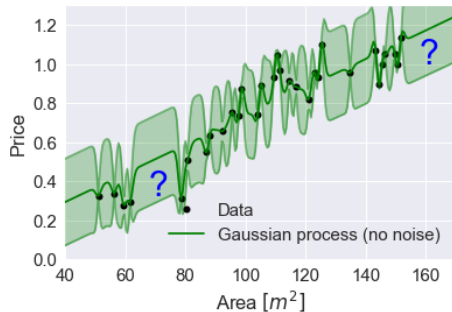
**Is the following statements true or false?**:

1. Gaussian processes can fit high non-linear functions, but the predictive means are given by a linear combination of the observations $\boldsymbol{y}$.

2. The variance of the posterior distribution is indepedent of the observations $\boldsymbol{y}$.

# What did we do?

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

$$p(f_*|\boldsymbol{f}) = \mathcal{N}\left(f_*\big|\boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{y}, K_{f_*f_*} - \boldsymbol{K}_{f_*f}\left(\boldsymbol{K}_{ff} + \sigma_{obs}^2\boldsymbol{I}\right)^{-1}\boldsymbol{K}_{f_*f}^T\right) \qquad (65)$$

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

# End of todays lecture

**Next lecture**:

- Kernels and covariance functions

- Model selection and hyperparameters

- Read ch. 4.2 and ch. 5.1-5.4 in Gaussian process book (gaussianprocess.org/gpml)

**Assignment**:

- Time to work on assignment #1 (deadline 20th of January)

- Should be handed in through the mycourses system

- In notebook format or in PDF with the same content