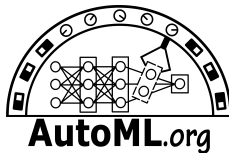


Automated Machine Learning (AutoML)

M. Lindauer F. Hutter

University of Freiburg



Lecture 6:

Bayesian linear regression and Gaussian processes



Where are we? The big picture

- Introduction
- Background
 - Design spaces in ML
 - Evaluation and visualization
- Hyperparameter optimization (HPO)
 - Bayesian optimization
 - Other black-box techniques
 - More details on Gaussian processes
- Pentecost (Holiday) – no lecture
- Architecture search I + II
- Meta-Learning
- Learning to learn & optimize
- Beyond AutoML: algorithm configuration and control
- Project announcement and closing



Today's Learning Goals

After today's lecture, you can ...

- Derive Bayesian linear regression
- Derive Gaussian processes
- Explain the relationship between Bayesian linear regression and Gaussian processes



Today's material is based on:

- Chapter 2 of
Gaussian Processes for Machine Learning, Carl Edward Rasmussen and Christopher K. I. Williams¹
- Another excellent resource: Philipp Hennig's Gaussian process tutorials²

¹<http://www.gaussianprocess.org/gpml/>

²<http://ei.is.tuebingen.mpg.de/person/phennig>

1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters

The univariate form

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x - \mu}{2\sigma^2}\right)$$

Named after Carl Friedrich Gauss (1777-1855)



Figure: Wikipedia, public domain

The univariate form

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

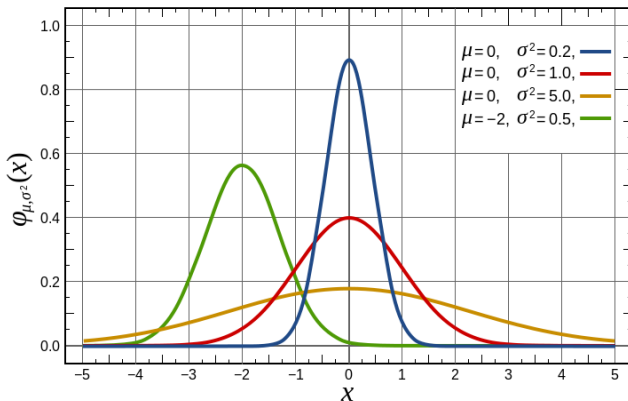


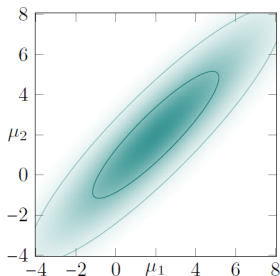
Figure: Wikipedia, public domain

The multivariate Gaussian (in N dimensions)

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

- \mathbf{x} and $\boldsymbol{\mu}$ are $N \times 1$ (column) vectors
- Σ is an $N \times N$ matrix

An example in 2 dimensions:



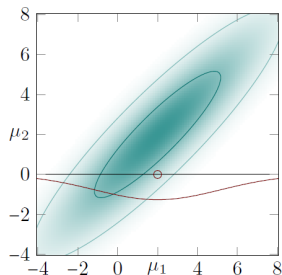
$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 \end{pmatrix}$$

Property 1: Closure under marginalization

Gaussian distributions are closed under marginalization

Let \mathbf{x}_1 and \mathbf{x}_2 be jointly Gaussian distributed:

$$\boldsymbol{\mu} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right) \text{ Then } \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}).$$



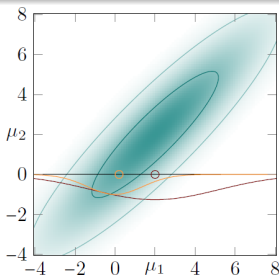
$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 \end{pmatrix}$$

Property 2: Closure under conditioning

Gaussian distributions are closed under conditioning

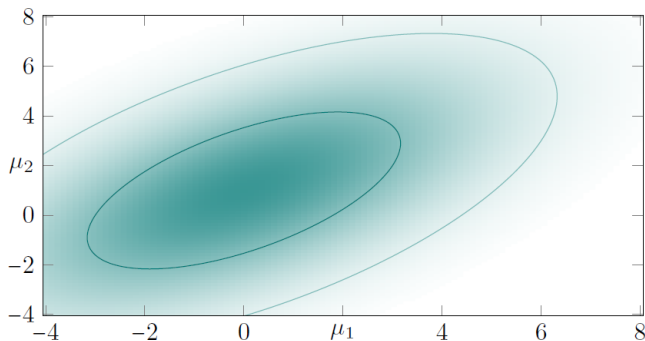
Let \mathbf{a} and \mathbf{b} be jointly Gaussian distributed: $\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \begin{bmatrix} A & C \\ C^\top & B \end{bmatrix} \right)$

Then $\mathbf{b}|\mathbf{a} \sim \mathcal{N}(\boldsymbol{\mu}_b + C^\top A^{-1}(\mathbf{a} - \boldsymbol{\mu}_a), B - C^\top A^{-1}C)$.



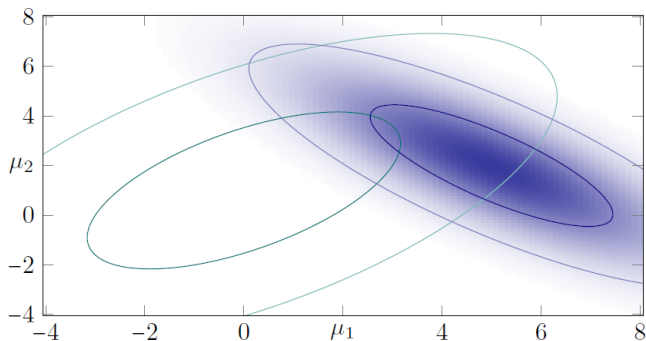
Property 3: Closure under multiplication

$$\mathcal{N}(x; a, A)$$



Property 3: Closure under multiplication

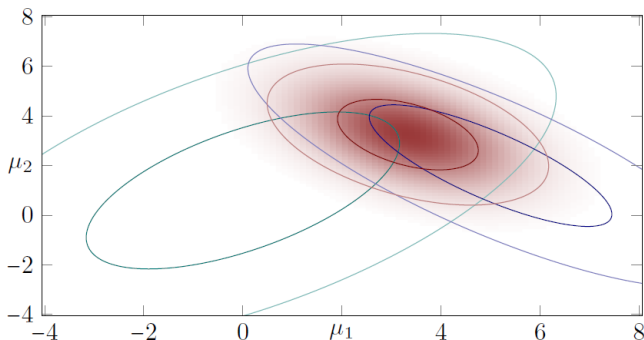
$$\mathcal{N}(x; b, B)$$



Property 3: Closure under multiplication

$$\mathcal{N}(x; a, A) \mathcal{N}(x; b, B) \propto \mathcal{N}(x; c, C)$$

$$C := (A^{-1} + B^{-1})^{-1} \quad c := C(A^{-1}a + B^{-1}b)$$



Why Gaussians?

- The central limit theorem tells us that means of large populations are distributed according to a Gaussian
- However, individual measurements are very rarely Gaussian-distributed
- Rather, Gaussians are mostly a mathematical convenience that lets us solve integrals in closed form!



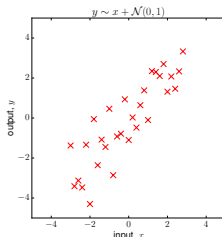
1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

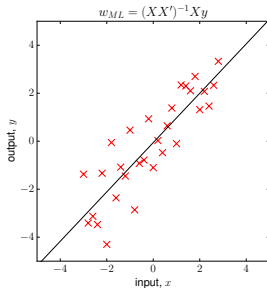
3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters



- We have a dataset of n i.i.d. observations, $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$
- Each \mathbf{x}_i denotes an input (column) vector of dimension D and each y_i is a scalar output or target
- We collect all inputs \mathbf{x}_i into the $D \times n$ design matrix $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]$ and the targets into the $n \times 1$ vector \mathbf{y} .
- Note: these slides are using the notation of Rasmussen's and Williams' book; \mathbf{X} is transposed from our usual use

Standard Linear Regression Model



$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \quad y = f(\mathbf{x}) + \epsilon \quad (2.1)$$

$$\epsilon \sim \mathcal{N}(0, \sigma_n^2) \quad (2.2)$$



The predictive distribution with fixed \mathbf{w}

$$\begin{aligned} p(\mathbf{y}|X, \mathbf{w}) &= \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^n \frac{1}{(\sqrt{2\pi}\sigma_n)} \exp\left(-\frac{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}|\mathbf{y} - \mathbf{X}^\top \mathbf{w}|^2\right) \\ &= \mathcal{N}(\mathbf{y}; \mathbf{X}^\top \mathbf{w}, \sigma_n^2 I) \end{aligned} \tag{2.3}$$



How does each of these steps follow?



The Bayesian way of setting w

- Define a **prior distribution over w** quantifying our uncertainty
- Integrate the likelihood $p(\mathbf{y}|\mathbf{X}, w)$ of the training data across all possible values of w , weighted by the prior $p(w)$

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, w)p(w)dw \quad (2.6)$$

- We don't commit to a single w , but integrate over an infinite number of possibilities

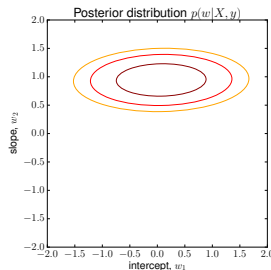
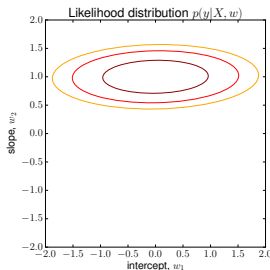
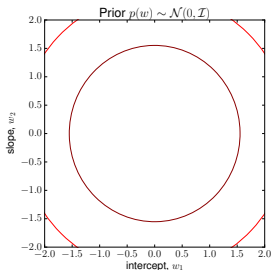


The Bayesian way of setting w

We combine the prior and the data likelihood using Bayes' rule:

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

$$p(w|y, X) = \frac{p(y|X, w)p(w)}{p(y|X)} \quad (2.5)$$



Posterior of \mathbf{w} under a Gaussian prior on \mathbf{w}

In general, the posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ has no closed-form solution, but it simplifies if we choose a prior $p(\mathbf{w})$ that is **conjugate** to the likelihood; in this case another Gaussian:

$$\mathbf{w} \sim \mathcal{N}(0, \Sigma_p) \quad (2.4)$$

In that case, the posterior is Gaussian again:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}\left(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} \mathbf{X} \mathbf{y}, A^{-1}\right), \quad (2.7)$$

with $A = \frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top + \Sigma_p^{-1}$.

We'll do the core of the derivation for this on the next two slides.



Key part of the derivation: completing the square (1)

- The data likelihood is $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \mathbf{X}^\top \mathbf{w}, \sigma_n^2 \mathbf{I})$.
- Written as a function of \mathbf{w} , this turns out to be proportional to a Gaussian $\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) = p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \propto \mathcal{N}(\mathbf{w}; \mathbf{w}', \Sigma^{-1})$
- We show this by matching coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w}$.



Key part of the derivation: completing the square (1)

- The data likelihood is $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \mathbf{X}^\top \mathbf{w}, \sigma_n^2 \mathbf{I})$.
- Written as a function of \mathbf{w} , this turns out to be proportional to a Gaussian $\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) = p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \propto \mathcal{N}(\mathbf{w}; \mathbf{w}', \Sigma^{-1})$
- We show this by matching coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w}$.

$$\begin{aligned}\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^\top (\mathbf{y} - \mathbf{X}^\top \mathbf{w})\right) \\ &\propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^\top (\mathbf{y} - \mathbf{X}^\top \mathbf{w})\right) \\ &= \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y} + (\mathbf{X}^\top \mathbf{w})^\top \mathbf{X}^\top \mathbf{w})\right)\end{aligned}$$



Key part of the derivation: completing the square (1)

- The data likelihood is $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \mathbf{X}^\top \mathbf{w}, \sigma_n^2 \mathbf{I})$.
- Written as a function of \mathbf{w} , this turns out to be proportional to a Gaussian $\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) = p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \propto \mathcal{N}(\mathbf{w}; \mathbf{w}', \Sigma^{-1})$
- We show this by matching coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w}$.

$$\begin{aligned}\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^\top (\mathbf{y} - \mathbf{X}^\top \mathbf{w})\right) \\ &\propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^\top (\mathbf{y} - \mathbf{X}^\top \mathbf{w})\right) \\ &= \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y} + (\mathbf{X}^\top \mathbf{w})^\top \mathbf{X}^\top \mathbf{w})\right)\end{aligned}$$

$$\begin{aligned}p(\mathbf{w}|\mathbf{w}', \Sigma) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}')^\top \Sigma^{-1}(\mathbf{w} - \mathbf{w}')\right) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{w}^\top \Sigma^{-1} \mathbf{w} - 2\mathbf{w}^\top \Sigma^{-1} \mathbf{w}' + \mathbf{w}'^\top \Sigma^{-1} \mathbf{w}')\right)\end{aligned}$$



Key part of the derivation: completing the square (2)

Let's match the coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w} \dots$

$$\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) \propto \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y} + (\mathbf{X}^\top \mathbf{w})^\top \mathbf{X}^\top \mathbf{w}) \right)$$

$$p(\mathbf{w} | \mathbf{w}', \Sigma) \propto \exp \left(-\frac{1}{2} (\mathbf{w}^\top \Sigma^{-1} \mathbf{w} - 2\mathbf{w}^\top \Sigma^{-1} \mathbf{w}' + \mathbf{w}'^\top \Sigma^{-1} \mathbf{w}') \right)$$



Key part of the derivation: completing the square (2)

Let's match the coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w} \dots$

$$\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) \propto \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y} + (\mathbf{X}^\top \mathbf{w})^\top \mathbf{X}^\top \mathbf{w}) \right)$$

$$p(\mathbf{w} | \mathbf{w}', \Sigma) \propto \exp \left(-\frac{1}{2} (\mathbf{w}^\top \Sigma^{-1} \mathbf{w} - 2\mathbf{w}^\top \Sigma^{-1} \mathbf{w}' + \mathbf{w}'^\top \Sigma^{-1} \mathbf{w}') \right)$$

This directly yields:

$$\begin{aligned} \mathbf{w}' &= \frac{1}{\sigma_n^2} \Sigma \mathbf{X} \mathbf{y} \\ \Sigma &= \left(\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top \right)^{-1} \end{aligned}$$



Key part of the derivation: completing the square (2)

Let's match the coefficients of \mathbf{w} and $\mathbf{w}^\top \mathbf{w} \dots$

$$\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) \propto \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y} + (\mathbf{X}^\top \mathbf{w})^\top \mathbf{X}^\top \mathbf{w}) \right)$$

$$p(\mathbf{w} | \mathbf{w}', \Sigma) \propto \exp \left(-\frac{1}{2} (\mathbf{w}^\top \Sigma^{-1} \mathbf{w} - 2\mathbf{w}^\top \Sigma^{-1} \mathbf{w}' + \mathbf{w}'^\top \Sigma^{-1} \mathbf{w}') \right)$$

This directly yields:

$$\begin{aligned} \mathbf{w}' &= \frac{1}{\sigma_n^2} \Sigma \mathbf{X} \mathbf{y} \\ \Sigma &= \left(\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top \right)^{-1} \end{aligned}$$

Thus: $\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) \propto \mathcal{N} \left(\frac{1}{\sigma_n^2} (\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{y}, (\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top)^{-1} \right)$.

Then, simply multiply two Gaussians to get the posterior in (2.7)



The posterior predictive distribution

Recall: in linear regression (so far without basis functions), each weight vector \mathbf{w} corresponds to a linear model:

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w} \quad (2.1)$$

We just derived the posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ over the weights \mathbf{w} :

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}\left(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} \mathbf{X} \mathbf{y}, A^{-1}\right). \quad (2.7)$$



The posterior predictive distribution

Recall: in linear regression (so far without basis functions), each weight vector \mathbf{w} corresponds to a linear model:

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w} \quad (2.1)$$

We just derived the posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ over the weights \mathbf{w} :

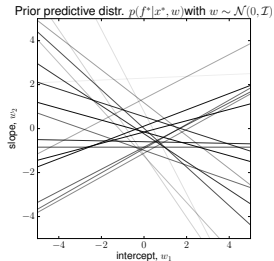
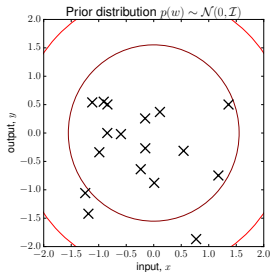
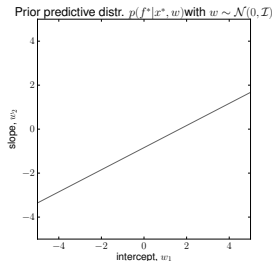
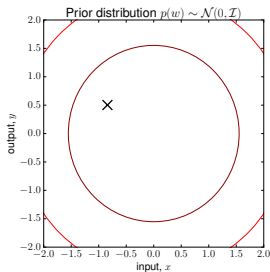
$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}\left(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} \mathbf{X} \mathbf{y}, A^{-1}\right). \quad (2.7)$$

Gaussians are closed under linear maps. Thus, the posterior predictive distribution at new input \mathbf{x}_\star given training data \mathbf{X}, \mathbf{y} is:

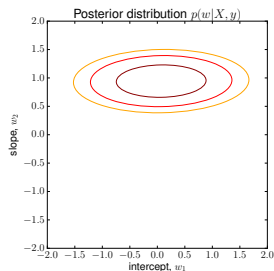
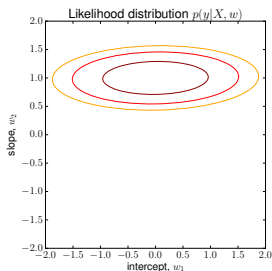
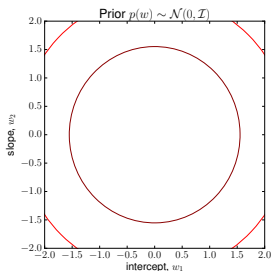
$$\begin{aligned} p(f_\star|\mathbf{x}_\star, \mathbf{X}, \mathbf{y}) &= p(\mathbf{x}_\star^\top \mathbf{w}|\mathbf{X}, \mathbf{y}) \\ &= \mathcal{N}(f_\star; \frac{1}{\sigma_n^2} \mathbf{x}_\star^\top A^{-1} \mathbf{X} \mathbf{y}, \mathbf{x}_\star^\top A^{-1} \mathbf{x}_\star) \end{aligned} \quad (2.9)$$



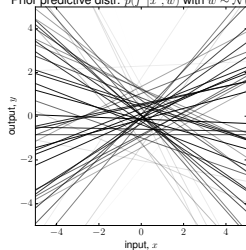
Weights correspond to functions (1)



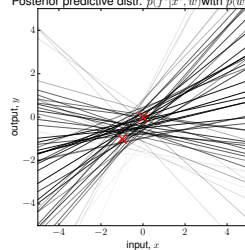
Weights correspond to functions (2)



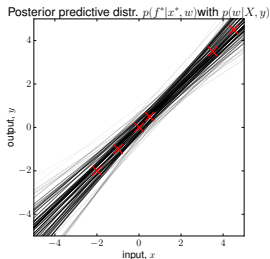
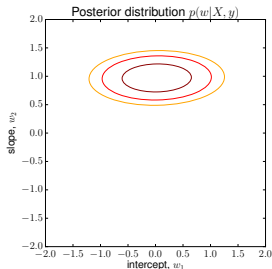
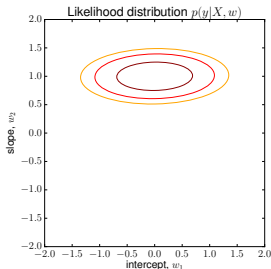
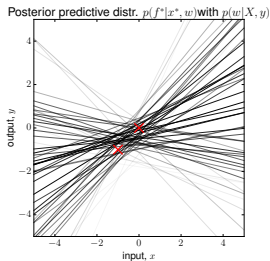
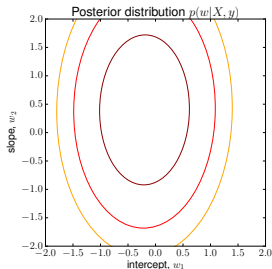
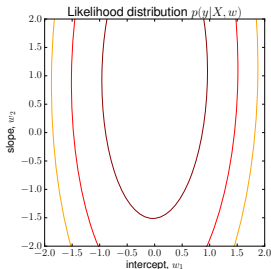
Prior predictive distr. $p(f^*[x^*, w])$ with $w \sim \mathcal{N}(0, \mathcal{I})$



Posterior predictive distr. $p(f^*[x^*, w])$ with $p(w|X, y)$



Enough data overwhelms the prior



1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters

Bayesian linear regression with basis functions

The same analysis goes through if we replace x with some so-called basis function of x , e.g., powers of x

$$\phi(x) = (1, x, x^2, x^3, \dots)^\top$$

Basis function linear regression is then:

$$f(x) = \phi(x)^\top \mathbf{w} \quad (2.10)$$

And the equivalent of 2.9 is:

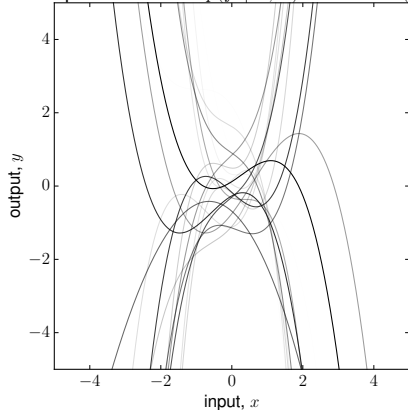
$$f_\star | \mathbf{x}_\star, \mathbf{X}, \mathbf{y} \sim \mathcal{N} \left(\frac{1}{\sigma_n^2} \phi(\mathbf{x}_\star)^\top A^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_\star)^\top A^{-1} \phi(\mathbf{x}_\star) \right) \quad (2.11)$$

with $\Phi = \phi(\mathbf{X})$ and $A = \sigma_n^{-2} \Phi \Phi^\top + \Sigma_p^{-1}$

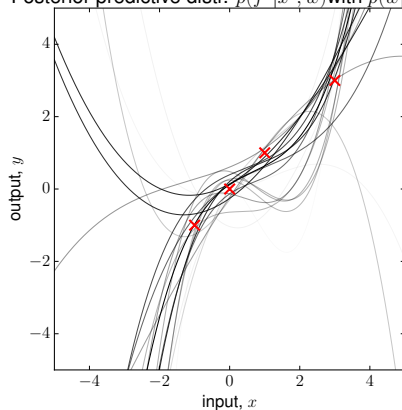


Bayesian Cubic Regression

Prior predictive distr. $p(f^*|x^*, w)$ with $w \sim \mathcal{N}(0, \mathcal{I})$



Posterior predictive distr. $p(f^*|x^*, w)$ with $p(w|X, y)$



Rewriting this equation using a kernel (1)

In (2.11), we have to invert the $N \times N$ matrix A
(where N = number of dimensions, or basis functions):

$$f_{\star} | \mathbf{x}_{\star}, \mathbf{X}, \mathbf{y} \sim \mathcal{N}\left(\frac{1}{\sigma_n^2} \phi(\mathbf{x}_{\star})^{\top} A^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_{\star})^{\top} A^{-1} \phi(\mathbf{x}_{\star})\right)$$

We can rewrite this equation to instead have to invert a $n \times n$ matrix
(where n = number of data points):

$$f_{\star} | \mathbf{x}_{\star}, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\phi_{\star}^{\top} \Sigma_p \Phi (\Phi^{\top} \Sigma_p \Phi + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.12) \\ \phi_{\star}^{\top} \Sigma_p \phi_{\star} - \phi_{\star}^{\top} \Sigma_p \Phi (\Phi^{\top} \Sigma_p \Phi + \sigma_n^2 I)^{-1} \Phi^{\top} \Sigma_p \phi_{\star}),$$

with shorthand notation ϕ_{\star} .

You will derive this equivalence as part of this week's exercise.



Rewriting this equation using a kernel (2)

$$f_{\star} | \mathbf{x}_{\star}, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\phi_{\star}^{\top} \Sigma_p \Phi (\Phi^{\top} \Sigma_p \Phi + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.12) \\ \phi_{\star}^{\top} \Sigma_p \phi_{\star} - \phi_{\star}^{\top} \Sigma_p \Phi (\Phi^{\top} \Sigma_p \Phi + \sigma_n^2 I)^{-1} \Phi^{\top} \Sigma_p \phi_{\star}),$$

with shorthand notation $\phi_{\star} = \phi(\mathbf{x}_{\star})$.

Note that the features $\phi(\mathbf{x}_{\star})$ only enter this equation in the forms of $\phi_{\star}^{\top} \Sigma_p \Phi$, $\phi_{\star}^{\top} \Sigma_p \phi_{\star}$, and $\Phi^{\top} \Sigma_p \Phi$.

We can replace all of these occurrences by expressions of the form

$$k(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^{\top} \Sigma_p \phi(\mathbf{x}_2)$$

This is known as the **kernel trick**. Then, the equation becomes:

$$f_{\star} | \mathbf{x}_{\star}, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(k_{\star}^{\top} (K + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.12) \\ k_{\star\star} - k_{\star}^{\top} (K + \sigma_n^2 I)^{-1} k_{\star}).$$



1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters

Definition (Gaussian process)

A **Gaussian process (GP)** is a collection of random variables, any finite number of which have a joint Gaussian distribution.

A GP is completely specified by its mean and covariance function:

$$\begin{aligned}m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})], \\k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]\end{aligned}$$

The GP is then:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (2.14)$$

The Bayesian linear regression model is a GP

The prior Bayesian linear model (before seeing any data \mathbf{X}, \mathbf{y}) is:

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}] = 0 \\ \mathbb{E}[f(\mathbf{x})f(\mathbf{x}')] &= \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi(\mathbf{x}') = \phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}')\end{aligned}\tag{2.15}$$

Thus $f(\mathbf{x})$ and $f(\mathbf{x}')$ are jointly Gaussian with mean zero and **covariance function** $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}')$



The squared exponential (SE) covariance

Many different covariance functions (or **kernel functions**) are possible; a prominent one is the squared exponential (SE) covariance function:

$$\text{cov}(f(\mathbf{x}_p), f(\mathbf{x}_q)) = k(\mathbf{x}_p, \mathbf{x}_q) = \exp\left(-\frac{1}{2}|\mathbf{x}_p - \mathbf{x}_q|^2\right) \quad (2.16)$$

- Points that are close to each other have correlation ≈ 1
- Points that are far away from each other have correlation ≈ 0
- This gives rise to very smooth functions

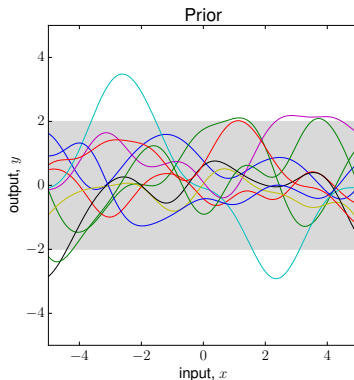


Samples from the GP prior

Samples for a finite number of points \mathbf{X}_\star

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

Thus, we're just sampling from a multivariate Gaussian distribution. For visualization, we connect the dots.



1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters

Prediction with Noise-free Observations

Observed function values f and new function values f_\star are jointly Gaussian-distributed

$$\begin{bmatrix} f \\ f_\star \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{X}, \mathbf{X}_\star) \\ K(\mathbf{X}_\star, \mathbf{X}) & K(\mathbf{X}_\star, \mathbf{X}_\star) \end{bmatrix} \right)$$

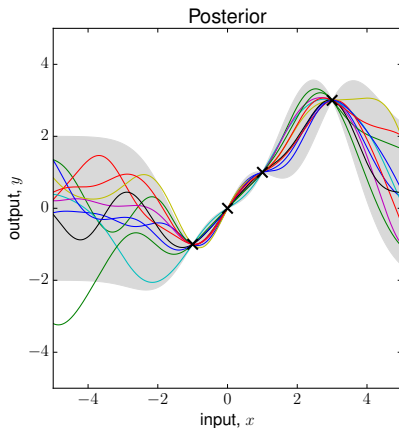
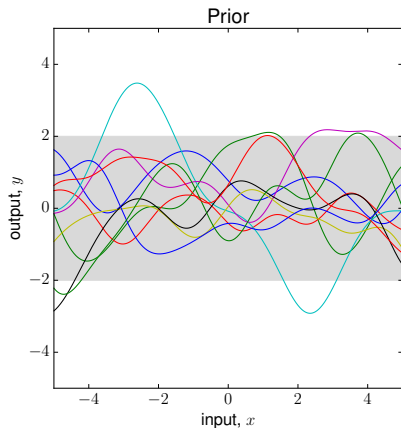
We simply apply the formula for conditioning a Gaussian to compute the predictive distribution $f_\star | f$:

$$f_\star | \mathbf{X}_\star, \mathbf{X}, f \sim \mathcal{N}(K(\mathbf{X}_\star, \mathbf{X})K(\mathbf{X}, \mathbf{X})^{-1}f, \quad (2.19) \\ K(\mathbf{X}_\star, \mathbf{X}_\star) - K(\mathbf{X}_\star, \mathbf{X})K(\mathbf{X}, \mathbf{X})^{-1}K(\mathbf{X}, \mathbf{X}_\star))$$



Samples from prior and posterior

Intuitively, we draw many (infinitely many) samples from the prior and discard those that don't agree perfectly with the noise-free data points



1 The Gaussian distribution

2 Gaussian processes: the weight space view

- Bayesian linear regression
- Bayesian linear regression with basis functions

3 Gaussian processes: the function space view

- The Case of Noise-free Observations
- The Case of Noisy Observations
- Marginal likelihood and kernel hyperparameters

Now we have noisy observations

As in the case of Bayesian linear regression:
we'll assume i.i.d Gaussian noise with variance σ_n^2

$$y = f(\mathbf{x}) + \epsilon$$
$$\epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

The prior on the noise observations is:

$$\begin{aligned} \text{cov}(y_p, y_q) &= k(\mathbf{x}_p, \mathbf{x}_q) + \sigma_n^2 \delta_{pq} \\ \text{cov}(\mathbf{y}) &= K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I \end{aligned} \tag{2.20}$$

Thus:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_\star \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I & K(\mathbf{X}, \mathbf{X}_\star) \\ K(\mathbf{X}_\star, \mathbf{X}) & K(\mathbf{X}_\star, \mathbf{X}_\star) \end{bmatrix} \right) \tag{2.21}$$



Predictions under noisy observations

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_\star \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I & K(\mathbf{X}, \mathbf{X}_\star) \\ K(\mathbf{X}_\star, \mathbf{X}) & K(\mathbf{X}_\star, \mathbf{X}_\star) \end{bmatrix} \right) \quad (2.21)$$

Using the same equation for conditioning Gaussians as before, we derive the posterior predictive distribution:

$$f_\star | \mathbf{X}, \mathbf{y}, \mathbf{X}_\star \sim \mathcal{N}(\bar{\mathbf{f}}_\star, \text{cov}(f_\star)), \text{ where} \quad (2.22)$$

$$\bar{\mathbf{f}}_\star \triangleq [f_\star | \mathbf{X}, \mathbf{y}, \mathbf{X}_\star] = K(\mathbf{X}_\star, \mathbf{X}) [K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I]^{-1} \mathbf{y}, \quad (2.23)$$

$$\text{cov}(f_\star) = K(\mathbf{X}_\star, \mathbf{X}_\star) - K(\mathbf{X}_\star, \mathbf{X}) [K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I]^{-1} K(\mathbf{X}, \mathbf{X}_\star) \quad (2.24)$$



The predictive distribution

Using more compact notation, we have:

$$\bar{f}_\star = \mathbf{k}_\star^\top (K + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.25)$$

$$\mathbb{V}[f_\star] = k(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{k}_\star^\top (K + \sigma_n^2 I)^{-1} \mathbf{k}_\star, \quad (2.26)$$

Note that with covariance function $\phi(\mathbf{x})^\top \Sigma_p \phi(\mathbf{x}')$ this is exactly the same as Bayesian linear regression:

$$f_\star | \mathbf{x}_\star, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\phi_\star^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.12) \\ \phi_\star^\top \Sigma_p \phi_\star - \phi_\star^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^\top \Sigma_p \phi_\star)$$



The predictive distribution

$$\bar{f}_\star = \mathbf{k}_\star^\top (K + \sigma_n^2 I)^{-1} \mathbf{y}, \quad (2.25)$$

$$\mathbb{V}[f_\star] = k(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{k}_\star^\top (K + \sigma_n^2 I)^{-1} \mathbf{k}_\star, \quad (2.26)$$

Some properties of the mean

- A linear combination of response values y
- A linear combination of n kernel functions, each one centered on a training data point:

$$f(\bar{\mathbf{x}}_\star) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_\star) \quad (2.27)$$

$$\text{where } \boldsymbol{\alpha} = (K + \sigma_n^2 I)^{-1} \mathbf{y}$$



- 1 The Gaussian distribution
- 2 Gaussian processes: the weight space view
 - Bayesian linear regression
 - Bayesian linear regression with basis functions
- 3 Gaussian processes: the function space view
 - The Case of Noise-free Observations
 - The Case of Noisy Observations
 - Marginal likelihood and kernel hyperparameters

The marginal likelihood

The marginal likelihood is the integral of the likelihood times the prior:

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f}. \quad (2.28)$$

It marginalizes (integrates out) the actual function values.
We know that this is distributed according to a Gaussian:

$$\mathbf{y} \sim \mathcal{N}(0, K + \sigma_n^2 I)$$

Usually, one works with the log marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}) = -\frac{1}{2}\mathbf{y}^\top (K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log |K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi \quad (2.30)$$

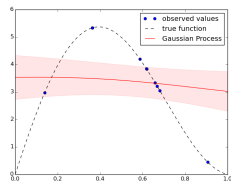
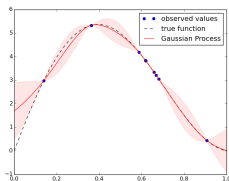
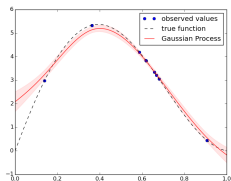


The effect of the kernel hyperparameters

The SE covariance function in 1-d:

$$k_y(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2l^2}(x_p - x_q)^2\right) + \sigma_n^2 \delta_{pq} \quad (2.31)$$

The marginal likelihood is highest for the left plot. It trades off data fit and “complexity” of the function



Summary by Learning Goals

Now (after digesting the material and doing the exercise), you can ...

- Derive Bayesian linear regression
- Derive Gaussian processes
- Explain the relationship between Bayesian linear regression and Gaussian processes

