

03 Bayesian Linear Regression

July 2, 2021

Information: *Basic concepts and simple examples of Bayesian linear regression*

Written by: *Zihao Xu*

Last update date: *07.02.2021*

1 Maximum Likelihood Estimation

1.1 Motivation

In the chapter talking about *Generalization and Regularization*, the concepts of parameter estimation, bias and variance are used to formally characterize notions of generalization, underfitting and overfitting. Here are some important remarks.

- View the parameter estimator $\hat{\theta}$ as a **function** of the sampled training dataset

$$\hat{\theta} = g(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$$

- The datasets (training, testing and probably validation) are generated by a **i.i.d.** probability distribution over datasets called the **data-generating process** (i.i.d. assumptions can be applied to almost all the common tasks)
- Assume that the true parameter value θ is fixed but unknown
- Since the **data** is drawn from a **random process**, any function of the data is random, which means the parameter estimator $\hat{\theta}$ is a **random variable**

The concepts of **bias** and **variance** are used to measure the performance of a parameter estimator. However, **for obtaining a good estimator**, it's not a good idea to guess that some function might make a good estimator and then to analyze its bias and variance. This motivated some principles from which specific functions that are good estimators for different models can be derived.

1.2 Definition

- Consider a set of m examples $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(\mathbf{x})$. Let $p_{\text{model}}(\mathbf{x}|\theta)$ be a parametric family of probability distributions over the same space indexed by θ
 - That is to say, p_{model} maps any configuration \mathbf{x} to a real number estimating the true probability $p_{\text{data}}(\mathbf{x})$

- Particularly, focus on the **likelihood** which is first introduced in the prerequisite chapter *Probability Theory*

$$\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \mid \boldsymbol{\theta} \sim p_{\text{model}}(\mathbf{x}^{(1:m)} \mid \boldsymbol{\theta})$$

As a fast review, the likelihood tells us how *plausible* it is to observe $\mathbf{x}^{(1:m)}$ if we know the model parameters are $\boldsymbol{\theta}$

- Since the examples are assumed to be drawn **independently**, the likelihood can be factorized

$$p_{\text{model}}(\mathbf{x}^{(1:m)} \mid \boldsymbol{\theta}) = \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)} \mid \boldsymbol{\theta})$$

Then **maximum likelihood** estimator for $\boldsymbol{\theta}$ is then defined as

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)} \mid \boldsymbol{\theta})$$

- While this simple production may lead to a lot of inconveniences such as **numerical underflow**, taking the **logarithm** of the likelihood does not change the location for maximum ($\arg \max$) but does conveniently transform a product into a sum

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)} \mid \boldsymbol{\theta})$$

- Obviously, rescaling the likelihood does not change the location for maximum ($\arg \max$), we can divide by m to obtain a version of the criterion that is expressed as an expectation with respect to the empirical distribution \hat{p}_{data} defined by the training data

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x} \mid \boldsymbol{\theta})]$$

- The most common choice for the likelihood of a single measurement is to pick it to be **Gaussian**

1.3 KL divergence

- Maximum likelihood estimation can be viewed as minimizing the dissimilarity between the empirical distribution \hat{p}_{data} , defined by the training set and the model distribution, with the degree of dissimilarity between the two measure by the **KL divergence**

$$D_{\text{KL}}(\hat{p}_{\text{data}} \parallel p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x} \mid \boldsymbol{\theta})]$$

The term on the left is a function only of the data-generating process, not the model. This means when we train the model to minimize the KL divergence, we need only minimize

$$-\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x} \mid \boldsymbol{\theta})]$$

- Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions. By definition, any loss consisting a negative log-likelihood is a **cross-entropy** between the **empirical distribution** defined by the training set (\hat{p}_{data}), and the **probability distribution** defined by the model (p_{model})

1.4 Conditional Log-Likelihood

- To apply *MLE* to most **supervised learning** tasks of predicting \mathbf{y} given \mathbf{x} , the maximum likelihood estimator is generalized to estimate a conditional probability $p_{\text{model}}(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$
- Consider a set of m examples $\mathcal{D} = \{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(\mathbf{x}, \mathbf{y})$. Factorize the data-generating process

$$p_{\text{data}}(\mathbf{x}, \mathbf{y}) = p_{\text{data}}(\mathbf{y} | \mathbf{x}) p_{\text{data}}(\mathbf{x})$$

Let $p_{\text{model}}(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})$ be a parametric family of probability distributions over the same space indexed by $\boldsymbol{\theta}$. It also can be factorized

$$p_{\text{model}}(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) = p_{\text{model}}(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) p_{\text{data}}(\mathbf{x})$$

Notice that the later part $p_{\text{data}}(\mathbf{x})$ is **fixed and shared**, the maximum likelihood estimation is going to focus on

$$p_{\text{model}}(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$$

Under the **i.i.d.** assumption, it can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

Similarly, this optimization problem is usually converted into a minimization problem by the **negative logarithm** operation considering computation issues

$$\boldsymbol{\theta}_{\text{ML}} = \arg \min_{\boldsymbol{\theta}} \left[- \sum_{i=1}^m \log \left[p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}) \right] \right]$$

1.5 Least Squares as Maximum Likelihood

- Least squares minimizing the mean square error is **equal** to maximum likelihood estimation when the likelihood is assigned to be **Gaussian**
- Assume the model is $\hat{y} = f(\mathbf{x}; \boldsymbol{\theta})$ with the dataset

$$\mathbf{X} = [\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots \quad \mathbf{x}^{(m)}], \mathbf{y} = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

The solution for $\boldsymbol{\theta}$ via least squares would be

$$\boldsymbol{\theta}_{\text{LS}} = \arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2$$

- From the point of view of maximum likelihood estimation, think of the model as producing a conditional distribution $p_{\text{model}}(y | \mathbf{x}, \boldsymbol{\theta})$ instead of producing a single prediction $\hat{y} = f(\mathbf{x}; \boldsymbol{\theta})$. Assign this likelihood of a single measurement to be Gaussian

$$p_{\text{model}}(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}, \sigma) = \mathcal{N}(y^{(i)} | f(\mathbf{x}; \boldsymbol{\theta}), \sigma^2)$$

where σ models the **noise**. This correspond to the belief that the measurement is around the model prediction $f(\mathbf{x}; \boldsymbol{\theta})$ but it is contained with Gaussian noise of variance σ^2 . For all the data, we have

$$\begin{aligned} p_{\text{model}}(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \sigma) &= \mathcal{N}(\mathbf{y} | f(\mathbf{X}; \boldsymbol{\theta}), \sigma^2 \mathbf{I}_m) \\ &= (2\pi)^{-m/2} \sigma^{-m} \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 \right) \end{aligned}$$

Then we have the maximum likelihood estimation to be

$$\begin{aligned}
\boldsymbol{\theta}_{\text{ML}} &= \arg \min_{\boldsymbol{\theta}} [-\log [p_{\text{model}}(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}, \sigma)]] \\
&= \arg \min_{\boldsymbol{\theta}} \left[\frac{m}{2} \log(2\pi) + m \log(\sigma) + \frac{1}{2\sigma^2} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 \right] \\
&= \arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 \\
&= \boldsymbol{\theta}_{\text{LS}}
\end{aligned}$$

Maximizing the likelihood with respect to $\boldsymbol{\theta}$ yields the same estimate as minimizing the squared error.

- The two criteria have **different values** but the **same location of the optimum**, which justifies the use of the LS as a maximum likelihood estimation procedure.
- Notice that σ is also a parameter to be optimized, maximize the likelihood with respect to σ

$$\begin{aligned}
\sigma_{\text{ML}} &= \arg \min_{\sigma} [-\log [p_{\text{model}}(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}, \sigma)]] \\
&= \arg \min_{\sigma} \left[\frac{m}{2} \log(2\pi) + m \log(\sigma) + \frac{1}{2\sigma^2} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 \right] \\
&= \arg \min_{\sigma} \left[m \log(\sigma) + \frac{1}{2\sigma^2} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 \right]
\end{aligned}$$

It can be easily solved by setting the derivative with respect to σ to zero

$$\begin{aligned}
m \frac{1}{\sigma_{\text{ML}}} - \frac{1}{\sigma_{\text{ML}}^3} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 &= 0 \\
m \sigma_{\text{ML}}^2 - \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2 &= 0 \\
\sigma_{\text{ML}}^2 &= \frac{1}{m} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2
\end{aligned}$$

- With the maximum likelihood estimation $\boldsymbol{\theta}_{\text{ML}}, \sigma_{\text{ML}}$, we can make **predictions** about y at a new point \mathbf{x}

$$p(y \mid \mathbf{x}, \boldsymbol{\theta}_{\text{ML}}, \sigma_{\text{ML}}) = \mathcal{N}(y \mid f(\mathbf{x}; \boldsymbol{\theta}_{\text{ML}}), \sigma_{\text{ML}}^2)$$

This means we're able to measure the **noise** occurred in sampling.

1.6 Properties of Maximum Likelihood Estimation

- Maximum likelihood estimator can be shown to be the **best** estimator asymptotically as the number of examples $m \rightarrow \infty$, in terms of its rate of convergence as m increases
- Whenever the **cross-entropy loss** or **least squares method** (e.g. MSE loss) is used, it can be viewed as computing the maximum likelihood estimate.
- Under appropriate conditions, the maximum likelihood estimator has the property of **consistency**
 - *Consistency: Unbiased*, converge to the true parameters as $m \rightarrow \infty$
 - The true distribution p_{data} must lie within the model family p_{model}
 - The true distribution p_{data} must correspond to exactly one value of $\boldsymbol{\theta}$

- The **statistical efficiency**, meaning that one consistent estimator may obtain lower generalization error for a fixed number of samples m , of the maximum likelihood estimator is very high among consistent estimators
 - It is asymptotically efficient, which means it achieves the Cramer Rao bound asymptotically
- It's mostly used when there is plenty of data
 - Tends to overfit when there is not enough data

2 Linear Regression via Maximum Likelihood Estimation

Here is a simple example showing how to apply **Maximum Likelihood Estimation** to linear regression and the correspondence between likelihood and least squares

2.1 Generate the dataset

The synthetic dataset is generated from

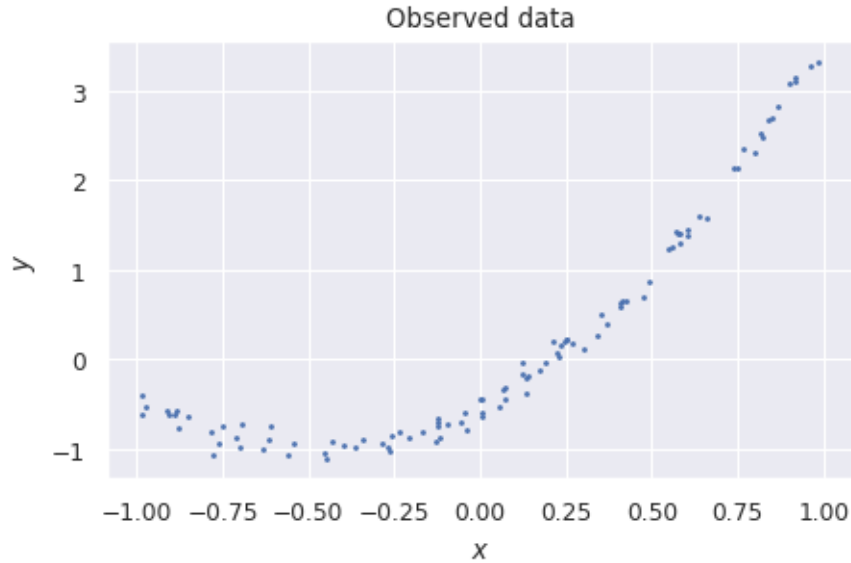
$$y_i = -0.5 + 2x_i + 2x_i^2 + 0.1\epsilon_i$$

where $\epsilon_i \sim \mathcal{N}(0, 1)$ and we sample $x_i \sim U([0, 1])$.

First generate this synthetic dataset and visualize the samples.

```
[1]: # Necessary modules
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
# Ensure reproducibility
np.random.seed(1234)
# Plot setting
sns.set()
sns.set_context('paper')

[2]: # Number of observations
num_obs = 100
# Sample x
x = (-1.0 + 2 * np.random.rand(num_obs)).reshape(-1, 1)
# True parameters
theta = np.array([-0.5, 2.0, 2.0]).reshape(-1, 1)
sigma = 0.1
# Calculate the corresponding y
y = theta[0] + theta[1] * x + theta[2] * x**2 \
    + sigma * np.random.randn(num_obs).reshape(-1, 1)
# Visualize the dataset
fig, ax = plt.subplots(figsize=(5, 3), dpi=100)
ax.plot(x, y, '.', markersize=2)
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')
ax.set_title('Observed data')
plt.show()
```



2.2 Build the model and train

As mentioned above, with setting the likelihood to be Gaussian, linear regression via least squares is a part of the procedure of maximum likelihood estimation.

$$\theta_{\text{ML}} = \theta_{\text{LS}}$$

Therefore, first get the values of the parameters via least squares. For convenience, here the solutions are got with the help of **Scikit-Learn**.

In this case, assume we know the exact degree of the polynomial so that we're not going to worry about validation and generalization issues, which makes us focusing on the MLE part. In the meantime, the dataset is not divided into training and validation dataset.

```
[3]: # Scikit-Learn Packages
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import PolynomialFeatures
# 'include_bias' and 'fit_intercept' cannot be both true
# Select one to be true according to personal preference
estimator = make_pipeline(PolynomialFeatures(2,include_bias=True),\
LinearRegression(fit_intercept=False))
# Train the model
estimator.fit(x,y)
# Compare the values
print("The true parameters are:\t",theta.T[0])
print("The estimated parameters are:\t",estimator[1].coef_[0])
```

The true parameters are: [-0.5 2. 2.]

The estimated parameters are: [-0.52407683 1.99641392 2.08372088]

2.3 Estimate the noise variance

In maximum likelihood estimation, one important thing is to also estimate the variance by maximizing the likelihood. Mention that the likelihood is set to be Gaussian, the variance should be

$$\sigma_{\text{ML}}^2 = \frac{1}{m} \|\mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta})\|_2^2$$

```
[4]: # Get the predictions of the trained model
y_pred = estimator.predict(x).reshape(-1,1)
# Calculate the estimated variance
sigma2_MLE = np.sum((y_pred-y)**2)/y.shape[0]
sigma_MLE = np.sqrt(sigma2_MLE)
# Compare the values
print("True sigma = %1.4f"%sigma)
print("MLE sigma = %1.4f"%sigma_MLE)
```

True sigma = 0.1000

MLE sigma = 0.0928

2.4 Make predictions

Now we have the maximum likelihood estimation $\boldsymbol{\theta}_{\text{ML}}, \sigma_{\text{ML}}$, we can make **predictions** about y at a new point \mathbf{x}

$$p(y|\mathbf{x}, \boldsymbol{\theta}_{\text{ML}}, \sigma_{\text{ML}}) = \mathcal{N}(y|f(\mathbf{x}; \boldsymbol{\theta}_{\text{ML}}), \sigma_{\text{ML}}^2)$$

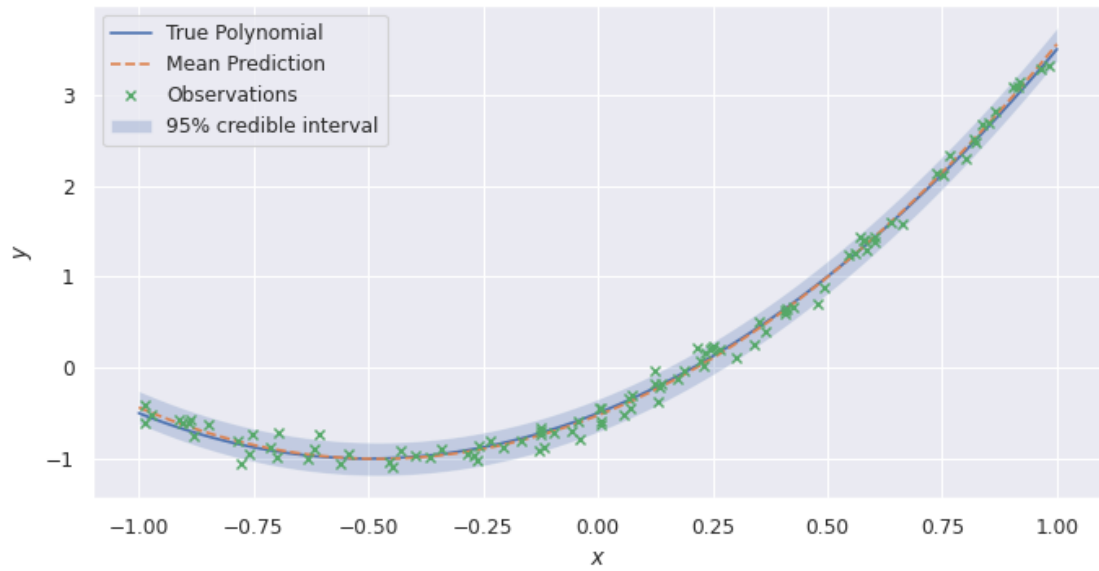
Let's visualize the 95% credible interval

$$(f(\mathbf{x}; \boldsymbol{\theta}_{\text{ML}}) - 1.96\sigma_{\text{ML}}, f(\mathbf{x}; \boldsymbol{\theta}_{\text{ML}}) + 1.96\sigma_{\text{ML}})$$

```
[5]: fig, ax = plt.subplots(figsize=(8,4),dpi=100)
# Points to be estimated
xe = np.linspace(-1,1,100)
# Predictions and true values
ye_pred = estimator.predict(xe.reshape(-1,1)).reshape(-1)
ye_true = theta[0] + theta[1] * xe + theta[2] * xe**2
# Lower bound and upper bound
ye_pred_lb = ye_pred - 1.96 * sigma_MLE
ye_pred_ub = ye_pred + 1.96 * sigma_MLE
# True polynomial
ax.plot(xe,ye_true,label='True Polynomial')
# Mean prediction
ax.plot(xe,ye_pred,'--',label='Mean Prediction')
# 95% credible interval
ax.fill_between(xe,ye_pred_lb,ye_pred_ub,alpha=0.25,label='95% credible_
→interval')
# Observations
ax.plot(x,y,'x',label='Observations')
# Labels
```



```
ax.set_xlabel('$x$')  
ax.set_ylabel('$y$')  
plt.legend(loc='upper left')  
plt.show()
```



3 Bayesian methods

3.1 Frequentist approach

- Review: In frequentist view, when we say that an outcome has a probability p of occurring, it means that if we repeated the experiment infinite times, then a proportion p of the repetitions would result in that outcome
- The true parameter value θ is supposed to be **fixed but unknown**
- The point estimate $\hat{\theta}$ is a random variable on account of it being a function of the dataset, which is seen as random
- Predictions are made based on estimating a single value of θ
- Addresses the uncertainty in a given point estimate of θ by evaluating its variance
 - The variance is an assessment of how the estimate might change with alternative samplings of the observed data
- **Maximum Likelihood Estimation** is a Frequentist approach

3.2 Bayesian Estimation

- Review: Bayesian uses the probability to reflect **degrees of certainty** in states of knowledge
- The true parameter θ is unknown or uncertain and thus is represented as a **random variable**. Before observing the data, represent our knowledge of θ using the **prior probability distribution** $p(\theta)$
 - Typically selects prior distribution that is quite broad (e.g. maximum entropy) to reflect a high degree of uncertainty in the value of θ before observing any data
- The dataset is directly observed and so is not random. Consider a set of data examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$, we can recover the effect of data on our belief about θ via Bayes' rule

$$p\left(\theta \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\right) = \frac{p\left(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \mid \theta\right) p(\theta)}{p\left(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\right)}$$

In the scenarios where Bayesian estimation is typically used, the prior begins as a relatively uniform or Gaussian distribution with high entropy, and the observation of the data usually causes the posterior to lose entropy and concentrate around a few highly likely values of the parameters

- Bayesian approach makes predictions using a full distance over θ . After observing m examples, the predicted distribution over the next data sample $\mathbf{x}^{(m+1)}$ is given by

$$p\left(\mathbf{x}^{(m+1)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\right) = \int p\left(\mathbf{x}^{(m+1)} \mid \theta\right) p\left(\theta \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\right)$$

Each value of θ with positive probability density contributes to the prediction of the next example, with the contribution weighted by the posterior density itself.

- Enable us to quantify the **epistemic uncertainty** induced by the limited number of observations used to estimate the weights

- The Bayesian prior distribution has an influence by shifting probability mass density towards regions of the parameter space that are preferred a priori.
 - In practice, the prior often expresses a **preference for models that are simpler** or more smooth.
 - This naturally arose **regularization** effect tends to protect well against overfitting.
- Bayesian methods typically generalize much better when limited training data is available but typically suffer from high computational cost when the number of training examples is large.
- Similarly, for **supervised learning** tasks, use the **conditional likelihood** instead.

3.3 Bayesian vs Frequentist Estimation

3.3.1 Frequentist Estimation

- Usually used for parameter estimation
- Typically low-bias high-variance estimates
- Most appropriate when the
 - prior information is weak
 - amount of data is large
 - quality of data is high

3.3.2 Bayesian Estimation

- Usually used for Machine Learning inference
- Typically high-bias low-variance estimates
- Most appropriate when the
 - prior information is strong
 - amount of data is small
 - quality of data is poor

3.4 Maximum A Posterior Estimation

- While most operations involving the Bayesian posterior for most interesting models are intractable, a point estimate offers a tractable approximation
- **Maximum A Posterior** point estimate benefits from the Bayesian approach by allowing the prior to influence the choice of the point estimate
- The **MAP** estimate chooses the point of maximal posterior probability (or maximal probability density in the more common case of continuous θ)

[]: