

Advanced Machine Learning

Cristian Perez Jensen

November 26, 2024

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	<i>Paradigms of data science</i>	1
2	<i>Anomaly detection</i>	2
3	<i>Density estimation</i>	6
4	<i>Regression</i>	9
4.1	<i>Linear regression</i>	9
4.2	<i>Polynomial regression</i>	10
4.3	<i>Regularization</i>	11
5	<i>Causality</i>	12
5.1	<i>Counterfactual invariance</i>	12
6	<i>Gaussian processes</i>	14
6.1	<i>Kernels</i>	15
6.2	<i>Prediction</i>	16

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . I.e., $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\boldsymbol{v} \in \mathbb{R}^n$	n -dimensional vector
$\boldsymbol{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
\boldsymbol{M}^\top	Transpose of matrix \boldsymbol{M}
\boldsymbol{M}^{-1}	Inverse of matrix \boldsymbol{M}
$\det(\boldsymbol{M})$	Determinant of \boldsymbol{M}
$\frac{\mathrm{d}}{\mathrm{d}x}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial \boldsymbol{x}}f(\boldsymbol{x})$	Partial derivative of $f(\boldsymbol{x})$ w.r.t. \boldsymbol{x} at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}f(\boldsymbol{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}^2f(\boldsymbol{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$

1 Paradigms of data science

Let $\{x_1, \dots, x_n\}$ be i.i.d. samples, generated by an unknown distribution P . Assume that this distribution is in a distribution family,

$$\mathcal{H} = \{p(\cdot \mid \theta) \mid \theta \in \Theta\}.$$

The goal is to learn the parameters θ that fit the data $\{x_1, \dots, x_n\}$ best.

Frequentism. In frequentism, the maximum likelihood estimator (MLE) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

Bayesianism. Bayesianism assumes that there is a prior over distributions. The maximum a posteriori (MAP) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\theta \mid \mathbf{X}) \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \cdot p(\theta) && \text{Bayes' rule.} \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\theta) + \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

In practice, the prior acts as a regularization term.

Statistical learning. Now, assume that we have labeled samples $\{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathcal{X} \times \mathcal{Y}$, where y is the target variable. Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. For a predictor function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we define its risk as the expected loss,

$$\mathcal{R}(f) \doteq \mathbb{E}_{X,Y}[\ell(y, f(x))].$$

In statistical learning, we want to find a function that minimizes the risk. However, since the distribution over X, Y is unknown, we cannot compute $\mathcal{R}(f)$ directly. Instead, we use the empirical risk as a surrogate,

$$\hat{\mathcal{R}}(f) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)).$$

The goal is to obtain the empirical risk minimizer,

$$f^* \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathcal{R}}(f),$$

where \mathcal{F} is a family of functions that we assume f belongs to.

2 Anomaly detection

In anomaly detection, we are given a sample of objects $\mathcal{X} \subseteq \mathbb{R}^d$ with a normal class $\mathcal{N} \in \mathcal{X}$ —the data points that we wish to keep. We wish to construct a function $\phi : \mathcal{X} \rightarrow \{0, 1\}$, such that

$$\phi(x) = 1 \iff x \notin \mathcal{N}.$$

Formally, an anomaly is an unlikely event. Hence, the strategy is to fit a model of a parametric family of distributions to the data \mathcal{X} ,

$$\mathcal{H} = \{p(\cdot | \theta) | \theta \in \Theta\}.$$

Then, we define the anomaly score of x as a low probability $p(x | \theta^*)$ according to the optimal model in this hypothesis class.

Anomaly detection in a high-dimensional space is hard, because the normal class can be very complex. The idea is to project \mathcal{X} down to a lower dimensionality and perform anomaly detection there—hopefully the projected version of the normal class $\Pi(\mathcal{N})$ is less complex. In order to find the optimal linear projection, we will use principal component analysis (PCA).

Furthermore, it has been observed that linear projections of high-dimensional distributions onto low-dimensional spaces resemble Gaussian distributions. Hence, after performing PCA, we will fit a Gaussian mixture model (GMM) to the projected data.

Principal component analysis. The goal of PCA is to linearly project \mathbb{R}^d to \mathbb{R}^{d^-} such that the maximum amount of variance of the data is preserved.¹ Consider the base case $d^- = 1$. Let $\mathbf{u} \in \mathbb{R}^d$ with $\|\mathbf{u}\| = 1$, we project onto \mathbf{u} by inner product,

$$x \mapsto \mathbf{u}^\top x.$$

The sample mean of the reduced dataset is computed by

$$\frac{1}{n} \sum_{i=1}^n \mathbf{u}^\top x_i = \mathbf{u}^\top \bar{x},$$

where \bar{x} is the sample mean of the original dataset. Further, the sample variance of the reduced dataset is

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\mathbf{u}^\top x_i - \mathbf{u}^\top \bar{x})^2 &= \frac{1}{n} \sum_{i=1}^n \mathbf{u}^\top (x_i - \bar{x})(x_i - \bar{x})^\top \mathbf{u} \\ &= \mathbf{u}^\top \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top \right) \mathbf{u} \\ &= \mathbf{u}^\top \Sigma \mathbf{u}, \end{aligned}$$

where Σ is the covariance matrix of the dataset. Since we want the projection that preserves the maximum variance, we have the following objective,

$$\mathbf{u}^* \in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u}.$$

¹ Components with larger variance are more informative.

The Lagrangian of this problem is

$$\mathcal{L}(\mathbf{u}; \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} + \lambda(1 - \|\mathbf{u}\|^2)$$

with gradient

$$\frac{\partial \mathcal{L}(\mathbf{u}; \lambda)}{\partial \mathbf{u}} = 2\Sigma \mathbf{u} - 2\lambda \mathbf{u} \stackrel{!}{=} 0.$$

So, \mathbf{u} must satisfy $\Sigma \mathbf{u} = \lambda \mathbf{u}$ — \mathbf{u} is an eigenvector of Σ . It is easy to see that this must be the principal eigenvector by rewriting the objective,

$$\begin{aligned} \mathbf{u}^* &\in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u} \\ &= \operatorname{argmax}_{\substack{\|\mathbf{u}\|=1 \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \|\mathbf{u}\|^2 \\ &= \operatorname{argmax}_{\substack{\mathbf{u} \in \mathbb{R}^d \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \\ &= \mathbf{u}_1. \end{aligned}$$

For $d^- > 1$, the remaining principal components can be computed with a similar idea. Iteratively, we factor out the previous principal components and do as above on the transformed dataset. For example, to get the second principal component, we first factor out the first principal component,

$$\mathcal{X}_1 \doteq \{\mathbf{x} - \operatorname{proj}_{\mathbf{u}_1}(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\} = \{\mathbf{x} - \mathbf{u}_1^\top \mathbf{x} \cdot \mathbf{u}_1 \mid \mathbf{x} \in \mathcal{X}\}.$$

Then, we do the same as above.

Gaussian mixture model. The probability density function (PDF) of a Gaussian mixture model with k components is formalized as a convex combination of Gaussians,

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \Sigma_j).$$

The parameters of this model are

$$\boldsymbol{\theta} = \{\pi_j, \boldsymbol{\mu}_j, \Sigma_j \mid j \in [k]\},$$

where $\sum_{j=1}^k \pi_j = 1$ and $\{\Sigma_j \mid j \in [k]\}$ are positive definite. We fit the parameters of this model by maximizing the log-likelihood,

$$\begin{aligned} \boldsymbol{\theta}^* &\in \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \log p(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}; \boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \log p(\mathbf{x}_i; \boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \log \sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \Sigma_j). \end{aligned}$$

```

1: Initialize  $\theta_0$ 
2: for  $t \in [T]$  do
3:    $q^* \in \operatorname{argmin}_q E(q, \theta_{t-1})$ 
4:    $\theta_t \in \operatorname{argmax}_\theta M(q^*, \theta)$ 
5: end for
6: return  $\theta_T$ 

```

Note that the above is intractable, so we would like to decompose it into tractable terms that can be computed. Let's assume that we know from which latent variable each data point was generated, then we can compute the MLE of the extended dataset $\{(x_i, z_i) \mid i \in [n]\}$ as

$$\begin{aligned}
\log p(\mathbf{X}, \mathbf{z}; \theta) &= \sum_{i=1}^n \log p(x_i, z_i) \\
&= \sum_{i=1}^n \log(\pi_{z_i} \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i})) \\
&= \sum_{i=1}^n \log \pi_{z_i} + \sum_{i=1}^n \log \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i}),
\end{aligned}$$

which is tractable to maximize. Let q be a distribution over $[k]$, then we can rewrite the log-likelihood into tractable terms,

$$\begin{aligned}
\log p(\mathbf{X}; \theta) &= \mathbb{E}_{z \sim q} [\log p(\mathbf{X}; \theta)] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(\mathbf{X}, z; \theta)}{p(z \mid \mathbf{X}; \theta)} \right) \right] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(\mathbf{X}, z; \theta)}{p(z \mid \mathbf{X}; \theta)} \frac{q(z)}{q(z)} \right) \right] \\
&= \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{p(\mathbf{X}, z; \theta)}{q(z)} \right]}_{\doteq M(q, \theta)} + \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid \mathbf{X}; \theta)} \right]}_{\doteq E(q, \theta)}.
\end{aligned}$$

These terms have the following two properties,

$$\begin{aligned}
\log p(\mathbf{X}; \theta) &\geq M(q, \theta), \quad \forall q, \theta \\
\log p(\mathbf{X}; \theta) &= M(q^*, \theta), \quad q^* = p(\cdot \mid \mathbf{X}; \theta), \quad \forall \theta.
\end{aligned}$$

Hence, we can use $M(q^*, \theta)$ as an approximation of $\log p(\mathbf{X}; \theta)$ around θ .

Theorem 2.1 (EM algorithm convergence). Using the expectation-maximization algorithm, $\{\log p(x; \theta_t)\}_{t=0}^T$ is non-decreasing.

Proof. Given a data point x and current parameters θ , we have the following update,

$$\theta' \in \operatorname{argmax}_{\theta \in \Theta} M(q^*, \theta).$$

Algorithm 1. The expectation-maximization algorithm, where

$$\begin{aligned}
M(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{p(\mathbf{X}, z; \theta)}{q(z)} \right] \\
E(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid \mathbf{X}; \theta)} \right].
\end{aligned}$$

Hence, we have

$$\log p(\mathbf{x}) = M(q^*, \boldsymbol{\theta}) \leq M(q^*, \boldsymbol{\theta}') \leq \log p(\mathbf{x}; \boldsymbol{\theta}').$$

Thus, $\{\log p(\mathbf{x}; \boldsymbol{\theta}_t)\}_{t=0}^T$ is non-decreasing. ■

Summary. In conclusion, given a set of data points \mathcal{X} with normal points $\mathcal{N} \subseteq \mathcal{X}$, we train an anomaly detector as follows,

1. Fit a projector $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^{d^-}$ using PCA;
2. Fit a probability density function $p(\cdot \mid \boldsymbol{\theta})$ with k components to $\{\pi(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\}$ using the EM algorithm;
3. For a point $\mathbf{x} \in \mathcal{X}$, its “anomaly score” is computed by $-\log p(\pi(\mathbf{x}); \boldsymbol{\theta})$.

3 Density estimation

In this section, we will consider parametric models,

$$\{p(\mathbf{x}; \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \Theta\}.$$

The problem we concern ourselves with is finding the best parameter $\boldsymbol{\theta}$. The most common method of finding the best parameters is maximum likelihood estimation (MLE),

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{MLE}} &\in \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \prod_{i=1}^n p(\mathbf{x}_i; \boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} - \sum_{i=1}^n \log p(\mathbf{x}_i; \boldsymbol{\theta}). \end{aligned}$$

The following properties makes the MLE estimator attractive,

1. Consistency—in the limit of n , $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ converges to the true parameter $\boldsymbol{\theta}^*$;
2. Equivariance—if $\hat{\boldsymbol{\theta}}$ is the MLE of $\boldsymbol{\theta}$, then $g(\hat{\boldsymbol{\theta}})$ is the MLE of $g(\boldsymbol{\theta})$ for any function g ;
3. Asymptotically normal—in the limit of n , $\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta} / \sqrt{n}$ converges to a random variable with distribution $\mathcal{N}(0, \mathcal{I}(\boldsymbol{\theta})^{-1})$, where \mathcal{I} is the Fisher information matrix;
4. Asymptotically efficient—in the limit of n , the MLE estimator has the smallest variance of all unbiased estimators.

We can understand the asymptotic efficiency (property 4) of estimators better by the Rao–Cramér bound, which provides a bound on the variance of the estimator. We will only consider the general scalar case, but it generalizes to the multivariate case.

Theorem 3.1 (Rao–Cramér bound (scalar case)). For any (scalar) unbiased estimator $\hat{\theta} : \mathcal{Y}^n \rightarrow \mathbb{R}$, given n data points, of $\theta \in \mathbb{R}$, we have the following bound on its variance,

$$\operatorname{Var}[\hat{\theta}(\mathbf{y})] \geq \frac{\left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1\right)^2}{\mathcal{I}_n(\theta)} + b_{\hat{\theta}}^2,$$

where $\mathcal{I}_n(\theta)$ is the Fisher information,

$$\mathcal{I}_n(\theta) \doteq \mathbb{E}_{\mathbf{y}|\theta} \left[\left(\frac{\partial}{\partial \theta} \log p(\mathbf{y} \mid \theta) \right)^2 \right] \stackrel{\text{iid}}{=} n \cdot \mathcal{I}(\theta),$$

and $b_{\hat{\theta}}$ is the bias of $\hat{\theta}$,

$$b_{\hat{\theta}} \doteq \mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}(\mathbf{y})] - \theta.$$

Proof. Let the “score” be defined as follows,

$$\Lambda(\mathbf{y}, \theta) \doteq \frac{\partial}{\partial \theta} \log p(\mathbf{y} | \theta) = \frac{1}{p(\mathbf{y} | \theta)} \frac{\partial}{\partial \theta} p(\mathbf{y} | \theta).$$

Then we have $\mathcal{I}_n(\theta) = \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta)^2]$.

The expected score is equal to zero,

$$\begin{aligned} \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta)] &= \int p(\mathbf{y} | \theta) \Lambda(\mathbf{y}, \theta) d\mathbf{y} \\ &= \int \frac{\partial}{\partial \theta} p(\mathbf{y} | \theta) d\mathbf{y} \\ &= \frac{\partial}{\partial \theta} \int p(\mathbf{y} | \theta) d\mathbf{y} \\ &= \frac{\partial}{\partial \theta} 1 \\ &= 0. \end{aligned}$$

(Hence, the Fisher information is equivalent to the variance of the score.)

Furthermore, the cross-correlation between $\Lambda(\mathbf{y}, \theta)$ and $\hat{\theta}(\mathbf{y})$ can be computed as

$$\begin{aligned} \text{Cov}_{\mathbf{y}|\theta} (\Lambda(\mathbf{y}, \theta), \hat{\theta}(\mathbf{y})) &= \mathbb{E}_{\mathbf{y}|\theta} [(\Lambda(\mathbf{y}, \theta) - \mathbb{E}[\Lambda(\mathbf{y}, \theta)]) (\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])] \\ &= \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y})] - \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta)] \mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}(\mathbf{y})] \\ &= \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y})] \\ &= \int p(\mathbf{y} | \theta) \Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\ &= \int \frac{\partial}{\partial \theta} p(\mathbf{y} | \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\ &= \frac{\partial}{\partial \theta} \int p(\mathbf{y} | \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\ &= \frac{\partial}{\partial \theta} (\mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}(\mathbf{y})] - \theta) + 1 \\ &= \frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1. \end{aligned}$$

Using the Cauchy-Schwarz inequality, we have the following bound,

$$\begin{aligned} \text{Cov}_{\mathbf{y}|\theta} (\Lambda(\mathbf{y}, \theta), \hat{\theta}(\mathbf{y}))^2 &= \left(\mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta) (\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])] \right)^2 \\ \left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1 \right)^2 &\leq \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta)^2] \cdot \mathbb{E}_{\mathbf{y}|\theta} [(\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])^2] \\ &= \mathcal{I}_n(\theta) \cdot \mathbb{E}_{\mathbf{y}|\theta} [(\hat{\theta}(\mathbf{y}) - \theta - \mathbb{E}[\hat{\theta}(\mathbf{y}) - \theta])^2] \\ &= \mathcal{I}_n(\theta) \cdot \left(\mathbb{E}_{\mathbf{y}|\theta} [(\hat{\theta}(\mathbf{y}) - \theta)^2] + (\mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}] - \theta)^2 \right. \\ &\quad \left. - 2\mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}(\mathbf{y}) - \theta] \right) \\ &\leq \mathcal{I}_n(\theta) \cdot \left(\mathbb{E}_{\mathbf{y}|\theta} [(\hat{\theta}(\mathbf{y}) - \theta)^2] - b_{\hat{\theta}}^2 \right). \end{aligned}$$

Re-arranging yields

$$\mathbb{E}_{\mathbf{y}|\theta} [(\hat{\theta}(\mathbf{y}) - \theta)^2] \geq \frac{\left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1 \right)^2}{\mathcal{I}_n(\theta)} + b_{\hat{\theta}}^2.$$

■

Note the trade-off between variance and bias. If $\frac{\partial}{\partial \theta} b_{\hat{\theta}} < 0$, then the variance might be smaller than the variance of the best unbiased estimator.

Corollary. Let $\hat{\theta}$ be an unbiased estimator—i.e., $b_{\hat{\theta}} = 0$ —then

$$\text{Var}[\hat{\theta}(\mathbf{y})] \geq \frac{1}{\mathcal{I}_n(\theta)}.$$

Lemma 3.2. The maximum likelihood estimator $\hat{\theta}_{\text{ML}}$ is asymptotically efficient,

$$\lim_{n \rightarrow \infty} \text{Var}[\hat{\theta}_{\text{ML}}] = \frac{1}{\mathcal{I}_n(\theta)}.$$

However, the MLE estimator is not necessarily efficient for finite samples.

4 Regression

In regression, we try to estimate a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that best fits a given dataset $\{(x_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^d \times \mathbb{R}$. I.e., we try to minimize the following loss function,

$$\ell(\beta) = \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|^2.$$

4.1 Linear regression

In linear regression, we assume the underlying data to be linear with a fixed noise,

$$Y | X = x \sim \mathcal{N}(\beta_*^\top x, \sigma^2)$$

for some ground truth β_* . We parametrize f as a linear function,

$$f(x; \beta) = \beta^\top x.$$

Under these assumptions, the minimizer of the loss function is analytically tractable—the ordinary least squares estimator (OLSE),

$$\hat{\beta} = (X^\top X)^{-1} X^\top y,$$

where $X \in \mathbb{R}^{d \times n}$ is the design matrix with respective outputs $y \in \mathbb{R}^n$.

In practice, it is important to remove outliers, because linear models can be heavily influenced by them. Also, the data should be standardized, such that all features are on the same scale, because differences in scale can make matrix inversion unstable.

TODO: Curse of dimensionality [Sur and Candès, 2019], Figures 4.1 and 4.2.

Further, if features are collinear, then the model might only learn the correlation with the target variable of one of them and discard the other. In addition, some singular values will equal zero. This makes matrix inversion unstable. This can easily be shown by considering the singular value decomposition $X = UDV^\top$, then

$$\begin{aligned} \hat{\beta} &= (X^\top X)^{-1} X^\top y \\ &= (VD^\top U^\top UDV^\top)^{-1} VD^\top U^\top y \\ &= (VD^\top DV^\top)^{-1} VD^\top U^\top y \\ &= V(D^{-1})^\top D^{-1} V^\top VD^\top U^\top y \\ &= V(D^{-1})^\top D^{-1} D^\top U^\top y \\ &= VD^{-1} U^\top y. \end{aligned}$$

The inversion of D is unstable if the singular values are small. The solution to this is to remove collinear features and/or to add regularization.

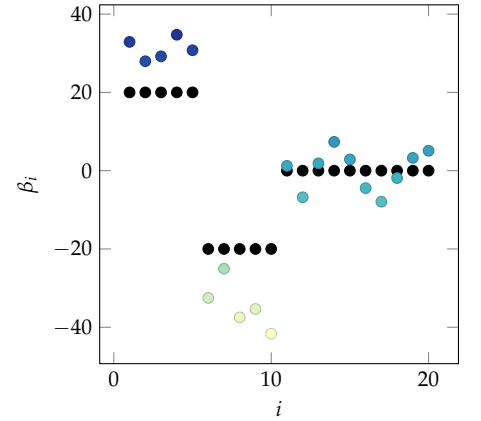


Figure 4.1. β^* is shown as black marks and $\hat{\beta}$ is indicated by the marks. As can be seen, $\hat{\beta}$ is overestimated for indices where $\beta_i^* \neq 0$.

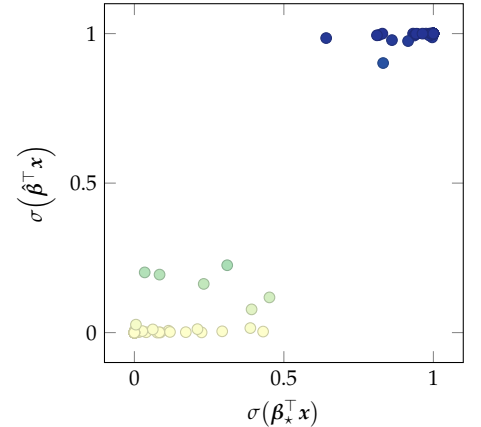


Figure 4.2. True *vs.* predicted probability—the estimator is overconfident in its predictions.

Consider the mean-squared error loss. The risk (expected loss) can be decomposed as follows into bias, variance, and noise terms,

$$\mathbb{E}[(f(\mathbf{x}) - y)^2] = \underbrace{(\mathbb{E}[f(\mathbf{x})] - \mathbb{E}[y])^2}_{\text{squared bias}} + \underbrace{\text{Var}[f(\mathbf{x})]}_{\text{variance}} + \underbrace{\mathbb{E}[(\mathbb{E}[y] - y)^2]}_{\text{noise}}.$$

The OLSE is the minimum variance unbiased estimator.² However, this does not mean it is the best, because introducing some bias may considerably decrease the variance. Bayesianism adds bias by introducing a prior—priors often induce a regularizing term.³

² This is proven by the Gauss-Markov theorem.

³ E.g., a Gaussian prior induces an ℓ_2 -norm regularizing term.

4.2 Polynomial regression

In polynomial regression, we construct a feature function of all possible polynomials, e.g.,

$$\phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, \dots].$$

We then perform linear regression in this space,

$$\psi(\mathbf{x}; \boldsymbol{\beta}) = \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle.$$

The problem is that this space is infinitely dimensional, and the inner product is ill-defined in an infinitely dimensional space.⁴ We can solve this by transforming the vector such that inner products cannot diverge by introducing a data-dependent scalar,

$$\phi(\mathbf{x}) = \alpha(\mathbf{x}) \cdot \tilde{\phi}(\mathbf{x}), \quad \alpha(\mathbf{x}) > 0.$$

⁴ This is due to $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \infty$ for some \mathbf{x}, \mathbf{x}' , e.g., $\mathbf{x} = \mathbf{x}' = [1, 1]$.

$\tilde{\phi}$ contains all polynomials in this equation.

An inner product w.r.t. this feature function is a valid inner product. Such transformations can have a closed form for the inner product, called kernelization. Commonly, the radial basis function (RBF) kernel is used,

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right), \quad \sigma \in \mathbb{R}.$$

Let $\Phi \in \mathbb{R}^{n \times \infty}$ contain all feature vectors, then the OLSE for polynomial regression can be computed by

$$\hat{\boldsymbol{\beta}} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}, \quad \Phi \in \mathbb{R}^{n \times \infty}.$$

However, $\Phi^\top \Phi$ cannot be computed, because it is $\infty \times \infty$ -dimensional. We can fix this by observing that $\Phi \Phi^\top \in \mathbb{R}^{n \times n}$ and rewriting the OLSE as

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} \\ &= (\Phi^\top \Phi)^{-1} \Phi^\top (\Phi \Phi^\top)^{-1} (\Phi \Phi^\top) \Phi^\top \mathbf{y} \\ &= (\Phi^\top \Phi)^{-1} \Phi^\top \Phi \Phi^\top (\Phi \Phi^\top)^{-1} \mathbf{y} \\ &= \Phi^\top (\Phi \Phi^\top)^{-1} \mathbf{y}. \end{aligned}$$

Let $\mathbf{K} = \Phi\Phi^\top$, then it only contains kernel evaluations— $k_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$. The next problem is that Φ^\top is still infinitely dimensional. However, this is not a problem, since when we want to make a prediction, we do the following,

$$\begin{aligned}\psi(\hat{\mathbf{x}}) &= \langle \phi(\hat{\mathbf{x}}), \hat{\boldsymbol{\beta}} \rangle \\ &= \phi(\hat{\mathbf{x}})^\top \Phi^\top \mathbf{K}^{-1} \mathbf{y}.\end{aligned}$$

Let $\mathbf{k}(\hat{\mathbf{x}}) = \phi(\hat{\mathbf{x}})^\top \Phi^\top$, then it only contains kernel evaluations with the new point— $k_i(\hat{\mathbf{x}}) = \langle \phi(\hat{\mathbf{x}}), \phi(\mathbf{x}_i) \rangle$. In conclusion, we can make predictions by

$$\psi(\hat{\mathbf{x}}) = \mathbf{k}(\hat{\mathbf{x}}) \mathbf{K}^{-1} \mathbf{y}, \quad \mathbf{k}(\hat{\mathbf{x}}) \in \mathbb{R}^{1 \times n}, \mathbf{K} \in \mathbb{R}^{n \times n}, \mathbf{y} \in \mathbb{R}^n.$$

However, the problem with this approach is that it takes $O(n^3)$ to make a prediction—it scales in the amount of data points.

4.3 Regularization

TODO: Ridge regression— ℓ_2 -norm.

TODO: LASSO— ℓ_1 -norm induces sparsity. Better interpretability properties.

TODO: Figure showing coefficient weights under different regularization factors.

5 Causality

In general, models do not distinguish between causal and non-causal factors in the feature space. Therefore, they might identify non-causal factors as highly correlating with the output variable. *E.g.*, when classifying images as either depicting a cow or a camel, the model might identify the background as an important feature, because cows are usually in grass and camels are usually in the desert. However, if this classifier were to see a cow with a different background, it would fail.

The following are examples of causal fallacies, where one might conclude that X causes Y ,

- Reverse causality, where Y actually predicts X and not the other way around;
- Third-cause fallacy, where there is a confounding factor Z that influences both X and Y ;
- Bidirectional causation, where X influence each other;

A domain shift happens when the test samples are drawn from a different distribution than the training samples. *E.g.*, COVID-19 detection models trained on a western population might not perform well on an eastern population.

Shortcut learning happens when there is a spurious correlation between causal and non-causal features in the training dataset. The resulting estimator abuses the non-causal features to generalize in the testing dataset. The solution to this is to encode the features such that they do not depend on the environment.

5.1 Counterfactual invariance

Let X be the feature vector representing the object and let Y be a target variable of interest, and let f be the function that estimates Y from X . Further, let W describe features that influence X , but should not influence the estimator f . Let a counterfactual be denoted as $X(w)$, which is the feature vector we would have obtained if we would have had $W = w$, after the fact. Then, the estimator f is counterfactually invariant if

$$f(X(w)) = f(X(w')), \quad \forall w, w' \in \text{range}(W).$$

In words, f is not influenced by the value of W . One way of obtaining a counterfactually invariant estimator is by extending the training dataset to contain enough counterfactuals. This can, for example, be achieved by data augmentation. However, this is not always possible.

Let X_A be the parts of X causally influenced by A , and let A^\perp be the set of variables independent of A . Then, f is counterfactually invariant if and only if f only depends on X_{W^\perp} .

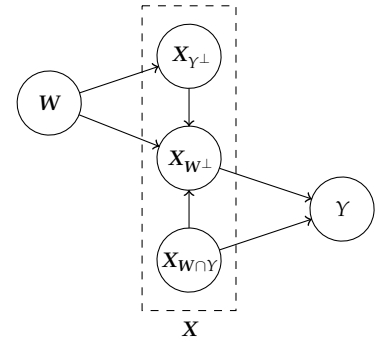


Figure 5.1. Causal graph.

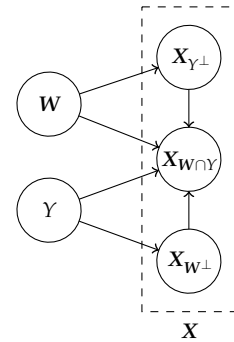


Figure 5.2. Anti-causal graph.

Theorem 5.1 (Necessary conditions for counterfactual invariance [Veitch et al., 2021]). If f is a counterfactually invariant predictor, then

- In the anti-causal scenario, $f(\mathbf{X}) \perp \mathbf{W} \mid Y$;
- In the causal scenario without selection (but possibly confounded), $f(\mathbf{X}) \perp \mathbf{W}$;
- In the causal scenario without confoundedness (but possibly with selection), $f(\mathbf{X}) \perp \mathbf{W} \mid Y$, as long as $\mathbf{X} \perp Y \mid \mathbf{X}_{\mathbf{W}^\perp}, \mathbf{W}$.

Proof. This can be proven by using d-separation. ■

6 Gaussian processes

Let the inputs be $\mathbf{X} \in \mathbb{R}^{n \times d}$, the outputs be $\mathbf{y} \in \mathbb{R}^n$, and the weight matrix be $\boldsymbol{\beta} \in \mathbb{R}^d$, then linear regression models the generative process as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

This is equivalent to defining a Gaussian over the predictions,

$$\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

BLR (*Bayesian Linear Regression*) extends linear regression by defining a prior over the regression coefficients,

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}^{-1}), \quad \boldsymbol{\Lambda} \in \mathbb{R}^{d \times d}.$$

The posterior can be analytically computed as

$$\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where

$$\boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}, \quad \boldsymbol{\Sigma} = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1}.$$

Proof.

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}) &\propto p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}) \cdot p(\boldsymbol{\beta}) && \text{Bayes' rule.} \\ &= \mathcal{N}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) \cdot \mathcal{N}(\boldsymbol{\beta}; \mathbf{0}, \boldsymbol{\Lambda}^{-1}) \\ &\propto \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} (\|\mathbf{y}\|^2 + \|\mathbf{X}\boldsymbol{\beta}\|^2 - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta}) + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) && \text{Cosine theorem.} \\ &= \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} + \frac{1}{\sigma^2} \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - \frac{2}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Lambda} \right) \boldsymbol{\beta} - \frac{2}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) - \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{X} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right)^\top \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) \right)\right) \\ &\propto \mathcal{N}\left(\boldsymbol{\beta}; \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}, \boldsymbol{\Sigma}\right). \end{aligned}$$

■

In general, we do not know the true weights that generated the data points. But, we can now define a joint distribution over output variables with unknown weights,

$$\mathbf{y} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{X} \boldsymbol{\Lambda}^{-1} \mathbf{X}^\top + \sigma^2 \mathbf{I}).$$

Proof. As we saw earlier, we compute outputs as follows,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

However, now the weights are unknown, so we make use of its prior to compute the distribution over \mathbf{y} ,

$$\begin{aligned} \mathbb{E}[\mathbf{y}] &= \mathbb{E}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}] \\ &= \mathbf{X}\mathbb{E}[\boldsymbol{\beta}] + \mathbb{E}[\boldsymbol{\epsilon}] \\ &= \mathbf{0}. \end{aligned}$$

$$\begin{aligned} \text{Cov}[\mathbf{y}] &= \mathbb{E}[(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})^\top] \\ &= \mathbb{E}[\mathbf{X}\boldsymbol{\beta}\boldsymbol{\beta}^\top \mathbf{X}^\top] + \mathbb{E}[\mathbf{X}\boldsymbol{\beta}\boldsymbol{\epsilon}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\beta}^\top \mathbf{X}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] \\ &= \mathbf{X}\mathbb{E}[\boldsymbol{\beta}\boldsymbol{\beta}^\top] \mathbf{X}^\top + \mathbf{X}\mathbb{E}[\boldsymbol{\beta}]\mathbb{E}[\boldsymbol{\epsilon}^\top] + \mathbb{E}[\boldsymbol{\epsilon}]\mathbb{E}[\boldsymbol{\beta}^\top] \mathbf{X}^\top + \sigma^2 \mathbf{I} \\ &= \mathbf{X}\boldsymbol{\Lambda}^{-1} \mathbf{X}^\top + \sigma^2 \mathbf{I}. \end{aligned}$$

$\boldsymbol{\beta}$ and $\boldsymbol{\epsilon}$ are independent.

■

GPs (*Gaussian Processes*) generalize BLR by observing that we can kernelize the covariance matrix,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \boldsymbol{\Lambda}^{-1} \mathbf{x}_j.$$

We could instead use any other kernel function to model other functions.

6.1 Kernels

Kernel functions specify the similarity between any two data points. They encode assumptions about the function that is to be learned.

Definition 6.1 (Kernel function). A kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ over $\mathcal{X} \subset \mathbb{R}^d$ must satisfy the following properties,

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= k(\mathbf{x}', \mathbf{x}) \\ \int k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' &\geq 0, \quad \forall f \in L_2. \end{aligned}$$

Symmetry.

Positive semi-definiteness.

Definition 6.2 (Stationary and isotropic). A kernel $k(\mathbf{x}, \mathbf{x}')$ is stationary if it only depends on $\mathbf{x} - \mathbf{x}'$. Further, it is isotropic if it only depends on $\|\mathbf{x} - \mathbf{x}'\|_2$.

The following are common kernels,

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

Linear kernel.

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + 1)^p, \quad p \in \mathbb{N}$$

Polynomial kernel.

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{\ell^2}\right), \quad \ell \in \mathbb{R}$$

RBF (*Radial Basis Function*) kernel.

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\kappa \mathbf{x}^\top \mathbf{x}') - b, \quad \kappa, b \in \mathbb{R}$$

Sigmoid kernel.

Different kernels have different invariance properties, such as invariance to rotation or translation. In order to learn invariances from data, many samples are needed. So, it might be better to encode them if we know them a priori.

Given two kernel functions k_1, k_2 defined on the same data space and positive scalar $c > 0$, the following are also valid kernels,

$$k(x, x') = k_1(x, x') + k_2(x, x')$$

$$k(x, x') = k_1(x, x') \cdot k_2(x, x')$$

$$k(x, x') = c \cdot k_1(x, x')$$

$$k(x, x') = \exp(k_1(x, x')).$$

In practice, the kernels are often composed together and hyperparameters are determined by performance on a held out validation dataset.

6.2 Prediction

As we saw earlier in the case of BLR, the marginal is jointly Gaussian,

$$\begin{bmatrix} y \\ y^* \end{bmatrix} \mid X, x^* \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K + \sigma^2 I & k(X, x^*) \\ k(x^*, X) & k(x^*, x^*) \end{bmatrix}\right).$$

Theorem 6.3 (Conditional Gaussian distribution). Given

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right).$$

Then, the conditional Gaussian is given by

$$x_2 \mid x_1 = z \sim \mathcal{N}\left(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(z - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right).$$

Using the above theorem, we can easily compute the predictive distribution,

$$y^* \mid x^*, X, y \sim \mathcal{N}\left(k^\top (K + \sigma^2 I)^{-1} y, c - k^\top (K + \sigma^2 I)^{-1} k\right),$$

where

$$K = K(X, X), \quad k = k(X, x^*), \quad c = k(x^*, x^*).$$

Now, we can compute a prediction along with its uncertainty. However, the problem with this approach is that it has $\mathcal{O}(n^3)$ runtime.

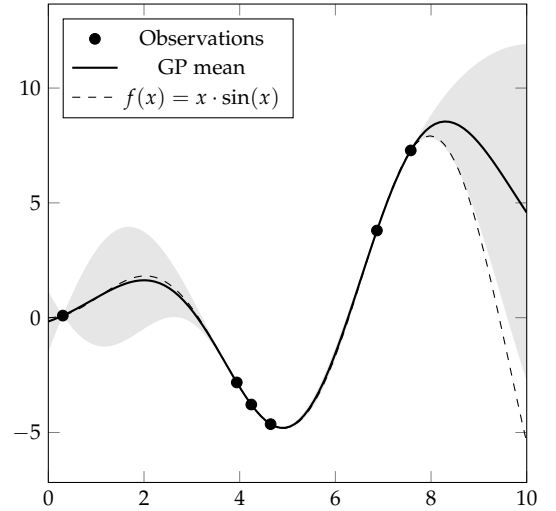


Figure 6.1. Fitted Gaussian process (RBF kernel, $\ell = 1.86$) with its 95% confidence interval.

References

- Pragya Sur and Emmanuel J Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.
- Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. *Advances in neural information processing systems*, 34:16196–16208, 2021.