

Machine Learning *Theory*

Christian Rossi

Academic Year 2023-2024

Abstract

The course topics are:

- Introduction: basic concepts.
- Learning theory:
 - Bias/variance tradeoff. Union and Chernoff/Hoeffding bounds.
 - VC dimension. Worst case (online) learning.
 - Practical advice on how to use learning algorithms.
- Supervised learning:
 - Supervised learning setup. LMS.
 - Logistic regression. Perceptron. Exponential family.
 - Kernel methods: Radial Basis Networks, Gaussian Processes, and Support Vector Machines.
 - Model selection and feature selection.
 - Ensemble methods: Bagging, boosting.
 - Evaluating and debugging learning algorithms.
- Reinforcement learning and control:
 - MDPs. Bellman equations.
 - Value iteration and policy iteration.
 - TD, SARSA, Q-learning.
 - Value function approximation.
 - Policy search. Reinforce. POMDPs.
 - Multi-Armed Bandit.

Contents

1	Introduction	1
1.1	Machine learning	1
1.1.1	Supervised learning	2
1.1.2	Unsupervised learning	2
1.1.3	Reinforcement learning	2
2	Supervised learning	4
2.1	Introduction	4
2.1.1	Function approximation	4
2.1.2	Taxonomy	5
2.2	Linear regression	5
2.2.1	Basis function	7
2.2.2	Regularization	8
2.2.3	Linear regression with probability	10
2.2.4	Challenges and limitations	13
2.3	Classification	13
2.3.1	Discriminant function	13
2.3.2	Probabilistic discriminative approaches	17
2.4	Kernel methods	18
2.4.1	Kernel function	19
2.4.2	Kernel ridge regression	21
2.4.3	Kernel regression	22
2.4.4	Gaussian processes	22
2.5	Support Vector Machines	23
2.5.1	Separable problems	23
2.5.2	Non-separable problems	25
2.5.3	Support vector machines training	25
2.5.4	Multi-class Support vector machines	26
3	Model evaluation	28
3.1	Bias-variance framework	28
3.1.1	Regularization and bias-variance	30
3.2	Model assessment	30
3.2.1	Optimal model	32
3.3	Model complexity	34
3.3.1	Feature selection	34
3.3.2	Dimensionality reduction	35

3.4	Ensemble	35
3.4.1	Bagging	35
3.4.2	Boosting	36
3.4.3	Summary	36
A	Algebra and statistics	37
A.1	Least squares	37
A.2	Matrices	38
A.2.1	Properties	38
A.3	Random variables	39
A.3.1	Continuous random variable	39
A.4	Distributions	40
A.5	Confidence intervals	41
A.6	hypothesis testing	41
A.6.1	Basic Gaussian test	41
A.6.2	P-value	42
A.7	Bayesian approach	42
B	Python introduction	43

CHAPTER 1

Introduction

1.1 Machine learning

Definition (*Learning*). A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , improves with experience E .

Machine learning, a subset of artificial intelligence, derives knowledge from experience and induction.

In machine learning, we depend on computers to make informed decisions using new, unfamiliar data. Designing a comprehensive set of meaningful rules can prove to be exceedingly difficult. Machine learning facilitates the automatic extraction of relevant insights from historical data and effectively applies them to new datasets.

The objective is to automate the programming process for computers, acknowledging the bottleneck presented by writing software. Instead, our aim is to utilize the data itself to accomplish the required tasks.

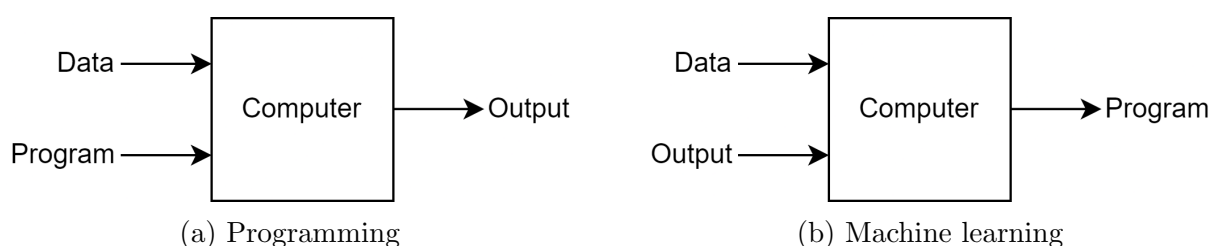


Figure 1.1: Difference between programming and machine learning

Machine learning paradigms can be categorized into three main types:

- *Supervised learning*: involves labeled data and direct feedback, aiming to predict outcomes or future events.
- *Unsupervised learning*: operates without labeled data or feedback, focusing on discovering hidden structures within the data.
- *Reinforcement learning*: centers around a decision-making process, incorporating a reward system to learn sequences of actions.

1.1.1 Supervised learning

Supervised learning encompasses several distinct tasks:

- *Classification*: this involves assigning predefined categories or labels to data points based on their features. The model is trained on labeled data, learning patterns to predict the class labels of new data points.
- *Regression*: the goal here is to predict continuous numerical values based on input features, as opposed to discrete class labels in classification. The model learns a function mapping input features to output values.
- *Probability estimation*: this task predicts the likelihood of certain events or outcomes occurring, often used to gauge the confidence of model predictions.

Formally, in supervised learning, a model learns from data to map known inputs to known outputs. The training set is denoted as $\mathcal{D} = \{\langle x, t \rangle\}$, Where $t = f(x)$, with f representing the unknown function to be determined using supervised learning techniques.

Various techniques can be employed for supervised learning, including linear models, artificial neural networks, support vector machines, and decision trees.

1.1.2 Unsupervised learning

Unsupervised learning encompasses two main tasks:

- *Clustering*: in this task, the objective is to group similar data points together based on their features, without predefined labels. The goal is to uncover underlying patterns or structures within the data. Clustering algorithms segment the data into clusters or groups, where data points within the same cluster exhibit greater similarity compared to those in different clusters. Unlike supervised learning, where labeled data is provided, clustering algorithms explore data solely based on features to identify similarities.
- *Dimensionality reduction*: this task involves reducing the number of input variables or features in a dataset while retaining essential information. This is often done to address the curse of dimensionality, enhance computational efficiency, and mitigate overfitting risks in models. Dimensionality reduction techniques aim to transform high-dimensional data into a lower-dimensional representation while preserving most relevant information.

Formally, in unsupervised learning, computers learn previously unknown patterns and efficient data representations. The training set is defined as $\mathcal{D} = \{x\}$, where the goal is to find a function f that extracts a representation or grouping of the data.

Various techniques are used for unsupervised learning, including k-means clustering, self-organizing maps, and principal component analysis.

1.1.3 Reinforcement learning

Reinforcement learning encompasses several key approaches:

- *Markov decision process*: a mathematical framework for modeling decision-making, involving states, actions, transition probabilities, and rewards. The goal is to find a policy that maximizes cumulative rewards while considering uncertainty.

- *Partially observable MDP*: an extension of MDP where the current state is uncertain and must be inferred from observations. The objective remains the same, but the agent maintains a belief over possible states based on observations.
- *Stochastic games*: models for decision-making with multiple agents, where outcomes depend on actions and random factors. Players aim to optimize strategies considering other players' actions and uncertainties.

In reinforcement learning, the computer learns the optimal policy based on a training set \mathcal{D} containing tuples $\langle x, u, x', r \rangle$, where x is the input, u is the action, x' is the resulting state after the action, and r is the reward. The policy Q^* is defined to maximize $Q^*(x, u)$ over actions u for each state x in the training set.

Various techniques such as Q-learning, SARSA, and fitted Q-iteration are used to find this optimal policy.

CHAPTER 2

Supervised learning

2.1 Introduction

Supervised learning stands as the predominant and well-established learning approach. Its core objective is to enable a computer, given a training set $\mathcal{D} = \{\langle x, t \rangle\}$, to approximate a function f that maps an input x to an output t . The input variables x , often referred to as features or attributes, are paired with output variables t , also known as targets or labels. The tasks undertaken in supervised learning are as follows:

- *Classification*: when t is discrete.
- *Regression*: when t is continuous.
- *Probability estimation*: when t represents a probability.

Supervised learning finds application in scenarios where:

- Humans lack the capability to perform the task directly (e.g., DNA analysis).
- Humans can perform the task but lack the ability to articulate the process (e.g., medical image analysis).
- The task is subject to temporal variations (e.g., stock price prediction).
- The task demands personalization (e.g., movie recommendation).

2.1.1 Function approximation

The process of approximating a function f from a dataset \mathcal{D} involves several steps:

1. *Define a loss function \mathcal{L}* : this function calculates the discrepancy between f and h , a chosen approximation.
2. *Select a hypothesis space \mathcal{H}* : this space consists of a set of candidate functions from which to choose an approximation h .
3. *Minimize \mathcal{L} within \mathcal{H}* : the goal is to find an approximation h within the hypothesis space \mathcal{H} that minimizes the loss function \mathcal{L} .

The hypothesis space \mathcal{H} can be expanded to theoretically achieve a perfect approximation of the function f . However, a significant challenge arises because the loss function \mathcal{L} cannot be easily determined, primarily due to the absence of the actual function f .

2.1.2 Taxonomy

The taxonomy is as follows:

- *Parametric* or *nonparametric*: parametric methods are characterized by having a fixed and finite number of parameters, while nonparametric methods have a number of parameters that depend on the training set.
- *Frequentist* or *Bayesian*: frequentist approaches utilize probabilities to model the sampling process, whereas Bayesian methods use probability to represent uncertainty about the estimate.
- *Empirical risk minimization* or *structural risk minimization*: empirical risk refers to the error over the training set, while structural risk involves balancing the training error with model complexity.

The type of machine learning can be:

- *Direct*: This method involves learning an approximation of f directly from the dataset \mathcal{D} .
- *Generative*: in this approach, the model focuses on modeling the conditional density $P(t|x)$ and then marginalizing to find the conditional mean:

$$\mathbb{E}[t|x] = \int t \cdot P(t|x) dt$$

- *Discriminative*: This method models the joint density $P(x, t)$, infers the conditional density $P(t|x)$, and then marginalizes to find the conditional mean:

$$\mathbb{E}[t|x] = \int t \cdot P(t|x) dt$$

2.2 Linear regression

The goal of regression is to approximate a function $f(x)$ that maps input x to a continuous output t from a dataset \mathcal{D} :

$$\mathcal{D} = \{\langle x, t \rangle\} \implies t = f(x)$$

Here, x is a vector. To perform regression, we assume the existence of a function capable of performing this mapping. The key components of constructing a linear regression problem include:

- The method used to model the function f (the hypothesis space).
- The evaluation criteria for the approximation (the loss function).
- The optimization process for optimizing the model.

In linear regression, the function f is modeled using linear functions. This choice is motivated by several factors:

- Linear models are easily interpretable, making them suitable for explanation.
- Linear regression problems can be solved analytically, allowing for efficient computation.
- Linear functions can be extended to model nonlinear relationships.
- More sophisticated methods often build upon or incorporate elements of linear regression.

Hypothesis space In mathematical terms, the approximation y can be defined as:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{D-1} w_j x_j = \mathbf{w}^T \mathbf{x}$$

Here, $\mathbf{x} = (1, x_1, \dots, x_{D-1})$ is a vector, and w_0 is called the bias parameter. It's important to note that the output y is a scalar value.

In a two-dimensional space, our hypothesis space will be the set of all points in the plane (w_0, w_1) . The coordinates of each point will correspond to a line in the (\mathbf{x}, y) space.

Loss function A commonly used error loss function for the linear regression problem is the sum of squared errors (SSE), defined as:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$$

This sum is also referred to as the residual sum of squares (RSS) and can be expressed as the sum of squared residual errors:

$$RSS(\mathbf{w}) = \|\epsilon_2^2\| = \sum_{i=1}^N \epsilon_i^2$$

This formulation of the loss function allows for obtaining a closed-form optimization solution.

Optimization For linear models, a closed-form optimization of the RSS, known as least squares, begins with the matrix representation of the loss function:

$$L(\mathbf{w}) = \frac{1}{2} RSS(\mathbf{w}) = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})$$

Here, $\Phi = [\phi(x_1) \ \dots \ \phi(x_N)]^T$ and $\mathbf{t} = [t_1 \ \dots \ t_N]^T$. To find the optimal \mathbf{w} , we compute the first derivative of $L(\mathbf{w})$ and set it to zero:

$$\hat{\mathbf{w}}_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

However, the inversion of the matrix $\Phi^T \Phi^{-1}$ can be computationally expensive, especially for large datasets, with a complexity of $O(nm^2 + m^3)$, assuming the matrix is non-singular (invertible).

To mitigate this, stochastic gradient descent (SGD) can be employed. The algorithm known as least mean squares (LMS) uses the following update rule:

$$L(\mathbf{x}) = \sum_n L(x_n)$$

Expanding this, we get:

$$\begin{aligned}\mathbf{w}^{(n+1)} &= \mathbf{w}^{(n)} - \alpha^{(n)} \nabla L(x_n) \\ &= \mathbf{w}^{(n)} - \alpha^{(n)} \left(\mathbf{w}^{(n)T} \phi(\mathbf{x}_n) - t_n \right) \phi(\mathbf{x}_n)\end{aligned}$$

Here, α is the learning rate, and convergence is guaranteed if $\sum_{n=0}^{\infty} \alpha^{(n)} = +\infty$ and $\sum_{n=0}^{\infty} \alpha^{(n)^2} < +\infty$.

If the regression problem involves multiple outputs, meaning that \mathbf{t} is not a scalar, we can solve each regression problem independently. However, we can still use the same set of basis functions. The solution for the weight vectors for all outputs can be expressed as:

$$\hat{\mathbf{W}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$$

Here, each column of matrix \mathbf{T} and $\hat{\mathbf{W}}$ corresponds to the target vector and the weight vector for each output, respectively. This solution can be easily decoupled for each output k :

$$\hat{\mathbf{w}}_k = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_k$$

An advantage of this approach is that $(\Phi^T \Phi)^{-1}$ only needs to be computed once, regardless of the number of outputs.

2.2.1 Basis function

While a linear combination of input variables may not always suffice to model data, we can still construct a regression model that is linear in its parameters. This can be achieved by defining a model using non-linear basis functions, expressed as:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Here, the components of the vector $\boldsymbol{\phi}(\mathbf{x}) = (1, \phi_1(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$ are referred to as features. These features allow for a more flexible representation of the input data, enabling the model to capture non-linear relationships between the input variables and the output.

Example:

Let's reconsider a set of data regarding individuals' weight and height, along with their completion times for a one-kilometer run:

Height (cm)	Weight (kg)	Completion time (s)
180	70	180
184	80	220
174	60	170

We can model this problem using a dummy variable and introduce the Body Mass Index (BMI) as a new feature:

Dummy variable	Height (cm)	Weight (kg)	BMI	Completion time (s)
x_0	x_1	x_2	x_3	t
1	180	70	21	180
1	184	80	23	220
1	174	60	20	170

Here, the dummy variable x_0 is always initialized to one. Now, we have the option to retain or discard the weight and height variables, considering only the BMI values for analysis.

The most commonly used basis functions in regression are:

- *Polynomial*:

$$\phi_j(x) = x^j$$

- *Gaussian*:

$$\phi_j(x) = \exp\left(-\frac{(x - \mu_j)^2}{2\sigma^2}\right)$$

- *Sigmoidal*:

$$\phi_j(x) = \frac{1}{1 + \exp\left(\frac{\mu_j - x}{\sigma}\right)}$$

Here, the constant μ_j is referred to as a hyperparameter, as its value needs to be determined through experimentation and depends on the user's experience.

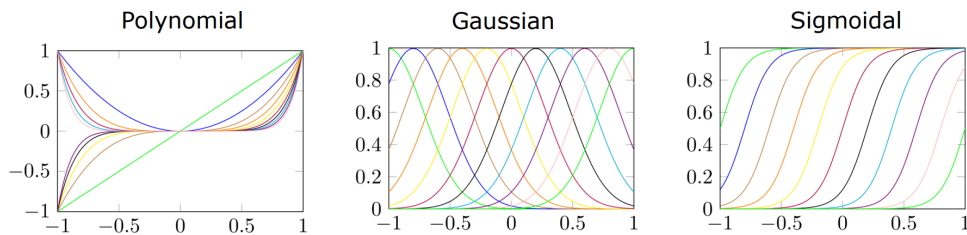


Figure 2.1: Some possible basis functions shapes

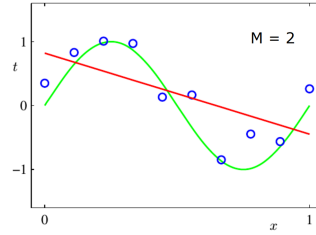
It's noteworthy that the Gaussian basis function allows for a local approximation by omitting values that are close to zero. This approach enables capturing the relationship between the input and output in a reduced input space area. As we move away from the mean, approaching zero, the values become negligible.

2.2.2 Regularization

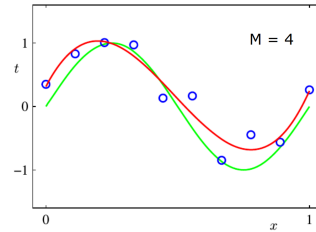
A function can achieve a better approximation by increasing the degree of the polynomial used in the regression.

Example:

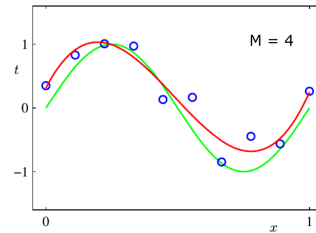
Consider a function generating a set of points with some noise:



Using a second-order polynomial instead of a linear one provides a better approximation:



Further improving the approximation can be achieved with a higher-degree polynomial (e.g., ninth grade):



However, increasing the polynomial degree also increases the complexity of the model parameters. To address this complexity, adjustments are needed in the loss function:

$$L(\mathbf{w}) = L_D(\mathbf{w}) + \lambda L_W(\mathbf{w})$$

Here, $L_D(\mathbf{w})$ represents the usual loss function, $L_W(\mathbf{w})$ reflects model complexity (a hyperparameter), and λ is the regularization coefficient. $L_W(\mathbf{w})$ can be tailored using ridge regression or lasso methods.

Ridge regression In ridge regression, the regularization term $L_W(\mathbf{w})$ is defined as:

$$L_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|_2^2$$

Thus, the overall loss function becomes:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(x_i))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Despite the regularization term, the loss function remains quadratic with respect to w , allowing for closed-form optimization:

$$\hat{\mathbf{w}}_{ridge} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

The term $\lambda \mathbf{I}$ is crucial in solving the singularity problem, as it transforms a non-singular matrix into a singular one with an appropriate choice of λ .

Lasso Another common regularization method is lasso, where the regularization term $L_W(\mathbf{w})$ is defined as:

$$L_W(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_1 = \frac{1}{2} \sum_{j=0}^{M-1} |w_j|$$

Thus, the overall loss function becomes:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(x_i))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_1$$

In this case, closed-form optimization is not possible. However, lasso typically leads to sparse regression models: when the regularization coefficient λ is large enough, some components of $\hat{\mathbf{w}}$ become equal to zero. Regularization can be seen as equivalent to minimizing $L_D(\mathbf{w})$ subject to the constraint:

$$\sum_{j=0}^{M-1} |w_j| \leq \eta$$

2.2.3 Linear regression with probability

We can approach regression in a probabilistic manner by defining a model that probabilistically maps inputs (x) to outputs (t). This model, denoted as $y(x, w)$, incorporates unknown parameters (w). We then model the likelihood, i.e., the probability that observed data \mathcal{D} is generated by a given set of parameters (w), as:

$$P(\mathcal{D}|\mathbf{w})$$

Finally, we estimate the parameters (w) by maximizing the likelihood:

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathcal{D}|\mathbf{w})$$

For linear regression, we define the model as:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon = \mathbf{w}^T \Phi(\mathbf{x}) + \epsilon$$

Here, we assume a linear model for $y(\mathbf{x}, \mathbf{w})$ and introduce noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Consequently, given a dataset \mathcal{D} of N samples with inputs $\mathbf{X} = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_N]$ and outputs $\mathbf{t} = [t_1 \ \cdots \ t_N]^T$, we have:

$$P(\mathcal{D}|\mathbf{w}) = P(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \Phi(\mathbf{x}_n), \sigma^2)$$

To find \mathbf{w}_{ML} , it is convenient to maximize the log-likelihood, obtaining:

$$\ell(\mathbf{w}) = \ln P(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} RSS(\mathbf{w})$$

Notice that the first part of the final formula is a constant independent of \mathbf{w} , so it can be ignored in maximizing the likelihood. Solving the optimization problem by setting the gradient to zero $\ell(\mathbf{w}) = 0$, yields the final formula:

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

This result aligns with the ordinary least squares approach.

This outcome allows us to interpret ordinary least squares from a probabilistic perspective, confirming that we are utilizing a normally distributed probabilistic function to generate the residuals in OLS.

Bayesian linear regression Bayesian linear regression follows a structured approach:

1. Formulation of probabilistic knowledge:
 - (a) Qualitatively define the model expressing our knowledge.
 - (b) Incorporate unknown parameters into the model.
 - (c) Represent assumptions about these parameters with a prior distribution before observing any data.

2. Data observation.

3. Computation of posterior probability distribution for parameters:

$$P(\text{parameters}|\text{data}) = \frac{P(\text{data}|\text{parameters})P(\text{parameters})}{P(\text{data})}$$

4. Utilization of Posterior Distribution to:

- Make predictions by averaging over the posterior distribution.
- Assess or accommodate uncertainty in parameter values.
- Make decisions by minimizing expected posterior loss.

The posterior distribution for model parameters is derived by combining the prior with the likelihood for parameters given the data:

$$P(\mathbf{w}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$$

Here, $P(\mathbf{w})$ represents the prior probability over parameters, $P(\mathcal{D}|\mathbf{w})$ denotes the likelihood, and $P(\mathcal{D})$ is the marginal likelihood acting as a normalization constant:

$$P(\mathcal{D}) = \int P(\mathcal{D}|\mathbf{w})P(\mathbf{w})d\mathbf{w}$$

The mode of the posterior, known as the Maximum A Posteriori (MAP) estimate, yields the most probable value of \mathbf{w} given the data.

A Gaussian likelihood assumption allows the prior to be modeled conveniently as a conjugate prior:

$$P(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{S}_0)$$

Consequently, the posterior remains Gaussian:

$$P(\mathbf{w}|\mathbf{t}, \Phi, \sigma^2) \propto \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{S}_0)\mathcal{N}(\mathbf{t}|\Phi\mathbf{w}, \sigma^2\mathbf{I})$$

Resulting in:

$$\begin{cases} P(\mathbf{w}|\mathbf{t}, \Phi, \sigma^2) = \mathcal{N}(\mathbf{w}|\mathbf{w}_N, \mathbf{S}_N) \\ \mathbf{w}_N = \mathbf{S}_N \left(\mathbf{S}_0^{-1}\mathbf{w}_0 + \frac{\Phi^T\mathbf{t}}{\sigma^2} \right) \\ \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \frac{\Phi^T\Phi}{\sigma^2} \end{cases}$$

Prior infinitely broad When the prior distribution is infinitely broad, the Maximum A Posteriori (MAP) estimate coincides with the Maximum Likelihood (ML) solution:

$$\begin{cases} \lim_{\mathbf{S}_0 \rightarrow \infty} \mathbf{w}_N = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \\ \lim_{\mathbf{S}_0 \rightarrow \infty} \mathbf{S}_N^{-1} = \frac{\Phi^T \Phi}{\sigma^2} \end{cases}$$

This yields the ordinary least squares formula. However, in this case, we also have the covariance matrix, providing information about the related uncertainty. The only missing parameter is σ^2 , which can be computed as:

$$\sigma^2 = \frac{1}{N - M} \sum_{n=1}^N (t_n - \hat{\mathbf{w}}^T(\phi)(\mathbf{x}_n))^2$$

The ML estimate \mathbf{w}_{ML} of \mathbf{w} has the smallest variance among linear unbiased estimates and the lowest Mean Squared Error (MSE) among linear unbiased estimates (Gauss-Markov).

Prior not infinitely broad When the prior distribution is not infinitely broad, such that $\mathbf{w}_0 = 0$ and $\mathbf{S}_0 = \tau^2 \mathbf{I}$, we can express the logarithm of the posterior distribution $P(\mathbf{w}|\mathbf{t})$ as:

$$\ln P(\mathbf{w}|\mathbf{t}) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 - \frac{1}{2\tau^2} \|\mathbf{w}\|_2^2$$

In this scenario, the Maximum A Posteriori estimate, MAP (\mathbf{w}_N), coincides with the solution of ridge regression $\hat{\mathbf{w}}_{ridge}$ with a regularization parameter λ set to $\lambda = \frac{\sigma^2}{\tau^2}$.

Sequential learning How to leverage the Bayesian approach for sequential learning:

1. Begin by computing the posterior with the initial data.
2. As additional data becomes available, update the prior with this new information to obtain the updated posterior.

Predictive distribution In a Bayesian framework, one can determine the probability distribution of the target variable for a new sample \mathbf{x}^* (given the training data \mathcal{D}) by integrating over the posterior distribution:

$$P(t^*|\mathbf{x}^*, \mathcal{D}) = \mathbb{E}[t^*|\mathbf{x}^*, \mathbf{w}, \mathcal{D}] = \int P(t^*|\mathbf{x}^*, \mathbf{w}, \mathcal{D}) P(\mathbf{w}|\mathcal{D}) d\mathbf{w}$$

This is commonly referred to as the predictive distribution. However, computing this predictive distribution typically involves the intractable task of determining the posterior distribution. Nevertheless, under certain assumptions, it is possible to compute the predictive distribution as follows:

$$\sigma_N^2(\mathbf{x}) = \sigma^2 + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

Here, as the number of data points N approaches infinity, the uncertainty associated with the parameters (second term) diminishes, and the variance of the predictive distribution depends solely on the variance of the data (σ^2).

2.2.4 Challenges and limitations

Modeling presents challenges in ensuring our model effectively represents a wide range of plausible functions while maintaining informative priors without overly spreading out probabilities or assigning negligible values.

On the computational side, limitations arise with analytical integration, particularly in cases involving non-conjugate priors and complex models. Approaches like Gaussian (Laplace) approximation, Monte Carlo integration, and variational approximation become necessary for addressing these complexities and achieving accurate results.

Linear models with fixed basis functions offer several benefits:

- They permit closed-form solutions, facilitating efficient computation.
- They lend themselves to tractable Bayesian treatment, enabling principled uncertainty quantification.
- They can capture non-linear relationships by employing appropriate basis functions.

However, these models also come with several drawbacks:

- Basis functions remain static and non-adaptive to variations in the training data.
- These models are susceptible to the curse of dimensionality, particularly when dealing with high-dimensional feature spaces.

2.3 Classification

Linear classification involves learning an approximation of a function $f(x)$ that maps inputs x to discrete classes C_k (with $k = 1, \dots, K$) from a dataset \mathcal{D} :

$$\mathcal{D} = \{\langle x, C_k \rangle\} \implies C_k = f(x)$$

Various approaches to classification include:

- *Discriminant function*: modeling a parametric function that directly maps inputs to classes and learning the parameters from the data.
- *Probabilistic discriminative approach*: designing a parametric model of $P(C_k|\mathbf{x})$ and learning the model parameters from the data.
- *Probabilistic generative approach*: modeling $P(\mathbf{x}|C_k)$ and class priors $P(C_k)$, fitting models to the data, and inferring the posterior using Bayes' rule:

$$P(C_k|\mathbf{x}) = \frac{P(\mathbf{x}|C_k)P(C_k)}{P(\mathbf{x})}$$

2.3.1 Discriminant function

We begin by examining the linear discriminant functions defined as:

$$f(\mathbf{x}, \mathbf{w}) = f\left(w_0 + \sum_{j=1}^{D-1} w_j x_j\right) = f(\mathbf{x}^T \mathbf{w} + w_0)$$

Here, the function $f(\cdot)$ is not linear in \mathbf{w} due to the presence of the (non-linear) activation function f , which yields either a discrete label or a probability value as its output.

The function $f(\cdot)$ partitions the input space into decision regions, with their boundaries known as decision boundaries or decision surfaces. Notably, these decision surfaces are linear functions of \mathbf{x} and \mathbf{w} , expressed as:

$$\mathbf{x}^T \mathbf{w} + w_0 = \text{constant}$$

It's important to note that generalized linear models are more complex to utilize compared to linear models due to the incorporation of non-linear activation functions.

Labels A common encoding for two-class problems involves binary encoding, where $t \in \{0, 1\}$. In this setup, $t = 1$ indicates the positive class, while $t = 0$ denotes the negative one. Using this encoding, both t and $f(\cdot)$ represent the probability of the positive class.

Another encoding option for two-class problems is $t \in \{-1, 1\}$, which is preferable for certain algorithms.

For problems with K classes, a typical choice is the 1-of- K encoding. Here, t is a vector of length K , with a 1 in the position corresponding to the encoded class. In this encoding scheme, both t and $f(\cdot)$ represent the probability density over the classes.

Example:

For instance, in a problem with $K = 5$ classes, a data sample belonging to class 4 would be encoded as:

$$t = [0 \ 0 \ 0 \ 1 \ 0]^T$$

Two-class problem The most general formulation for a discriminant linear function in a two-class linear problem is:

$$f(\mathbf{x}, \mathbf{w}) = \begin{cases} C_1 & \text{if } \mathbf{x}^T \mathbf{w} + w_0 \geq 0 \\ C_2 & \text{otherwise} \end{cases}$$

From this formulation, we can deduce the following properties:

- The decision boundary is $y(\cdot) = \mathbf{x}^T \mathbf{w} + w_0 = 0$.
- The decision boundary is orthogonal to \mathbf{w} .
- The distance of the decision boundary from the origin is $\frac{w_0}{\|\mathbf{w}\|_2}$.
- The distance of the decision boundary from \mathbf{x} is $\frac{y(\mathbf{x})}{\|\mathbf{w}\|_2}$.

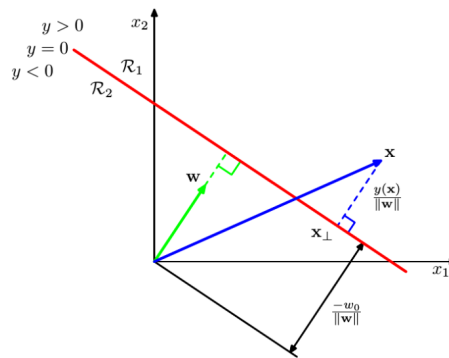


Figure 2.2: Two-class decision problem boundaries

Multiple-class problem In multiple class problems with K classes, various encoding methods can be employed:

- *One versus the rest*: this approach involves using $K - 1$ binary classifiers, where each classifier distinguishes between one class (C_i) and the rest of the classes. However, this method introduces ambiguity since there may be regions mapped to multiple classes.
- *One versus one*: this method utilizes $\frac{K(K-1)}{2}$ binary classifiers, where each classifier discriminates between pairs of classes C_i and C_j . Similar to the one versus the rest approach, this method also suffers from ambiguity.

One potential solution to mitigate the ambiguity in multi-class classification is to employ K linear discriminant functions:

$$y_k(\mathbf{x}) = \mathbf{x}^T \mathbf{w}_k + w_{k0} \quad k = 1, \dots, K$$

In this approach, an input vector \mathbf{x} is assigned to class C_k if $y_k > y_j$ for all $j \neq k$. This method ensures that the decision boundaries are singly connected and convex.

Linear basis function models Up to this point, we have focused on models operating within the input space. However, we can enhance these models by incorporating a fixed set of basis functions $\phi(\mathbf{x})$. Essentially, this involves applying a non-linear transformation to map the input space into a feature space. Consequently, decision boundaries that are linear within the feature space would correspond to nonlinear boundaries within the input space. This extension enables the application of linear classification models to problems where samples are not linearly separable.

Ordinary least squares Let's consider a K -class problem using a 1-of- K encoding for the target. Each class is modeled with a linear function:

$$y_k(\mathbf{x}) = \mathbf{x}^T \mathbf{w}_k + w_{k0} \quad k = 1, \dots, K$$

In matrix notation, this can be expressed as:

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}$$

Here, $\tilde{\mathbf{W}}$ is of size $(D + 1) \times K$, where its k -th column is denoted by $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$, and $\tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$.

Given a dataset $\mathcal{D} = \{\mathbf{x}_i, \mathbf{t}_i\}$ where $i = 1, \dots, N$, we can utilize the least squares method to determine the optimal value of $\tilde{\mathbf{W}}$, resulting in:

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{T}}$$

Here, $\tilde{\mathbf{X}}$ is an $N \times (D + 1)$ matrix with its i -th row being $\tilde{\mathbf{x}}_i^T$ and $\tilde{\mathbf{T}}$ is an $N \times K$ matrix with its i -th row as \mathbf{t}_i^T . In this setup, any new sample $\tilde{\mathbf{x}}_{new}^T$ is assigned to class C_k if $t_k > t_j$ for all j , where t_k represents the k -th component of the model output computed as $t_k = \tilde{\mathbf{x}}^T \tilde{\mathbf{w}}_k$.

The primary challenge with employing ordinary least squares in classification is that the resulting decision boundaries between regions can vary significantly based on the distribution of the data. This method may yield effective or suboptimal boundaries depending on the characteristics of the dataset.

Perceptron To address the issue of poor boundaries, one approach is to utilize a model known as the perceptron. Proposed by Rosenblatt in 1958, the perceptron is a linear discriminant model designed specifically for two-class problems, with class encoding as $\{-1, 1\}$. The perceptron model is defined as:

$$f(\mathbf{x}, \mathbf{w}) = \begin{cases} +1 & \text{if } \mathbf{x}^T \mathbf{w} + w_0 \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

The perceptron algorithm aims to determine a decision surface, also known as a separating hyperplane, by minimizing the distance of misclassified samples to the boundary. This minimization of the loss function can be achieved using stochastic gradient descent.

Although simpler loss functions could theoretically be used, they are often more complex to minimize in practice. Therefore, stochastic gradient descent is commonly employed for optimization in perceptron learning.

The core concept of the perceptron is to optimize \mathbf{w} such that $\mathbf{w}^T \phi(\mathbf{x}_i) \geq 0$ for $\mathbf{x}_i \in C_1$ and $\mathbf{w}^T \phi(\mathbf{x}_i) < 0$ otherwise. The perceptron criterion is expressed as:

$$L_P \mathbf{w} = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n$$

Here, correctly classified samples do not contribute to L , and each misclassified sample $\mathbf{x}_i \in \mathcal{M}$ contributes as $\mathbf{w}^T \phi(\mathbf{x}_n) t_n$.

Minimizing L_P is achieved using stochastic gradient descent:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha \nabla L_P(\mathbf{w}) = \mathbf{w}^{(k)} + \alpha \phi(\mathbf{x}_n) t_n$$

Since the scale of \mathbf{w} does not affect the perceptron function, the learning rate α is often set to 1. The perceptron algorithm takes a dataset $\mathcal{D} = \{\mathbf{x}_i, \mathbf{t}_i\}$ where $i = 1, \dots, N$.

Algorithm 1 Perceptron algorithm

```

1: Initialize  $\mathbf{w}_0$ 
2:  $k \leftarrow 0$ 
3: repeat
4:    $k \leftarrow k + 1$ 
5:    $n \leftarrow k \bmod N$ 
6:   if then  $\hat{t}_n \neq t_n$ 
7:      $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + \phi(\mathbf{x}_n) t_n$ 
8:   end if
9: until convergence

```

Theorem 2.3.1 (Perceptron convergence). *If the training dataset is linearly separable in the feature space Φ , then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps.*

Several steps may be necessary, making it challenging to distinguish between non-separable problems and slowly converging ones. If multiple solutions exist, the one obtained by the algorithm depends on the parameter initialization and the order of updates.

2.3.2 Probabilistic discriminative approaches

In a discriminative approach, we model the conditioned class probability directly:

$$P(C_1|\phi) = \frac{1}{1 + e^{-\mathbf{w}^T \phi}} = \sigma(\mathbf{w}^T \phi)$$

Here, $\sigma(\cdot)$ denotes the sigmoidal function. This model is commonly referred to as logistic regression.

Maximum likelihood Given a dataset $\mathcal{D} = \{\mathbf{x}_i, t_i\}$, where $i = 1, \dots, N$ and $t_i \in \{0, 1\}$, we aim to maximize the likelihood, i.e., the probability of observing the targets given the inputs $P(\mathbf{t}|\mathbf{X}, \mathbf{w})$. We model the likelihood of a single sample using a Bernoulli distribution, employing the logistic regression model for conditioned class probability:

$$P(t_n|\mathbf{x}_n, \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1-t_n} \quad y_n = P(t_n = 1|\mathbf{x}_n, \mathbf{w}) = \sigma(\mathbf{w}^T \phi_n)$$

Assuming independent sampling of data in \mathcal{D} , we have:

$$P(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{(1-t_n)} \quad y_n = \sigma(\mathbf{w}^T \phi_n)$$

The negative log-likelihood (also known as cross-entropy error function) serves as a convenient loss function to minimize:

$$L(\mathbf{w}) = -\ln P(\mathbf{t}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)) = \sum_{n=1}^N L_n$$

The derivative of L yields the gradient of the loss function:

$$\nabla L(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

Due to the nonlinearity of the logistic regression function, a closed-form solution is not feasible. Nevertheless, the error function is convex, allowing for gradient-based optimization (even in an online learning setting).

Multi class logistic regression In multi class problems, $P(C_k|\phi)$ is modeled by applying a softmax transformation to the output of K linear functions (one for each class):

$$P(C_k|\phi) = y_k(\phi) = \frac{e^{\mathbf{w}_k^T \phi}}{\sum_j e^{\mathbf{w}_j^T \phi}}$$

Similar to the two-class logistic regression and assuming 1-of- K encoding for the target, we compute the likelihood as:

$$P(\mathbf{T}|\Phi, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \left(\prod_{k=1}^K P(C_k|\phi_n)^{t_{nk}} \right) = \prod_{n=1}^N \left(\prod_{k=1}^K y_{nk}^{t_{nk}} \right)$$

As in the two-class problem, we minimize the cross-entropy error function:

$$L(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln P(\mathbf{T}|\Phi, \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \left(\sum_{k=1}^K t_{nk} \ln y_{nk} \right)$$

Then, we compute the gradient for each weight vector:

$$\nabla L_{\mathbf{w}_j}(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

Perceptron Replacing the logistic function with a step function in logistic regression yields the same updating rule as the perceptron algorithm.

2.4 Kernel methods

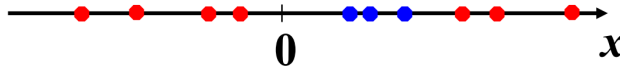
Frequently, we seek to detect nonlinear patterns within our datasets. In nonlinear regression, the connection between input and output may deviate from linearity, while in nonlinear classification, class boundaries might not be linearly separable. Linear models often prove insufficient in capturing such complexities. However, kernel methods offer a solution by transforming data into higher-dimensional spaces where linear relationships become apparent, thereby enabling linear models to effectively operate in nonlinear scenarios.

The process of transforming the original input space into a feature space is termed feature mapping, denoted as:

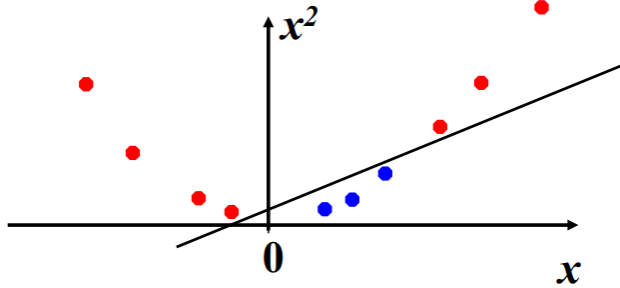
$$\Phi : x \rightarrow \phi(x)$$

Example:

Consider a binary classification problem where no linear separator exists:



Now, let's map the input space (a single variable x) to a feature space with two features: $x \rightarrow \{x, x^2\}$. As a result, the data becomes linearly separable:



This concept extends naturally to higher dimensions and more intricate problem domains.

However, a significant drawback arises known as the curse of dimensionality. This occurs due to the exponential growth in the number of features as the input variables increase, rendering the mapping computationally infeasible. Kernel methods offer a solution to this challenge

by bypassing the need for explicit computation of the feature mapping. While they are computationally intensive, they remain feasible for practical implementation.

2.4.1 Kernel function

The kernel function is defined as the scalar product between the feature vectors of two data samples:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

The kernel function exhibits symmetry: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$. It can be interpreted as a measure of similarity between \mathbf{x} and \mathbf{x}' .

Interestingly, very large feature vectors, even infinite ones, can result in a kernel function that is computationally tractable.

Certain special classes of kernel functions exist:

- *Stationary kernels*: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$.
- *Homogeneous kernels* (or radial basis functions): $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$.

Kernel function design We are not obligated to compute the kernel function by first generating the feature space, as we aim to avoid explicitly calculating the feature vectors. Two primary approaches exist for designing a kernel function:

- Create kernel functions directly from scratch.
- Design kernel functions by applying a predefined set of rules to existing ones.

In both cases, it's crucial to ensure that the resulting kernel functions are valid, meaning they correspond to a scalar product in some feature space.

Theorem 2.4.1 (Mercer). *Any continuous, symmetric, positive semi-definite kernel function $k(\mathbf{x}, \mathbf{x}')$ can be expressed as a dot product in a high-dimensional space.*

For this theorem, the necessary and sufficient condition for a function $k(\mathbf{x}, \mathbf{x}')$ to be a valid kernel is that the Gram matrix \mathbf{K} is positive semi-definite for all possible choices of $\mathcal{D} = \{\mathbf{x}_i\}$. This condition implies that $\mathbf{x}^T \mathbf{K} \mathbf{x} > 0$ for any non-zero real vector \mathbf{x} , meaning that the double sum $\sum_i \sum_j \mathbf{K}_{ij} \mathbf{x}_i \mathbf{x}_j$ is strictly positive for any real numbers \mathbf{x}_i and \mathbf{x}_j .

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ the following rules can be applied to design a new valid kernel:

1. $k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$, where $c > 0$ is a constant.
2. $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$, where $f(\cdot)$ is any function.
3. $k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$, where $q(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(\mathbf{x}, \mathbf{x}') = e^{k_1(\mathbf{x}, \mathbf{x}')}.$
5. $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$.
6. $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$.
7. $k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$, where $\phi(\mathbf{x})$ maps \mathbf{x} to \mathbb{R}^M and $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M .

8. $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where \mathbf{A} is a symmetric semidefinite matrix.
9. $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$.
10. $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$.

Kernel trick It's feasible to modify the representation of linear models by substituting terms involving $\phi(\mathbf{x})$ with alternatives solely based on $k(\mathbf{x}, \cdot)$. In essence, the output of a linear model can be computed solely based on the similarities between data samples, as computed with the kernel function.

This methodology, known as the kernel trick, finds application in various learning algorithms including: ridge regression, $K - NN$ regression, perceptron, nonlinear PCA, and support vector machines.

Gaussian kernel The Gaussian kernel is a widely employed kernel function in various machine learning algorithms. Its mathematical representation is given by:

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}}$$

This kernel function defines a similarity measure between two vectors \mathbf{x} and \mathbf{x}' in the feature space. It assigns higher similarity to vectors that are closer to each other, based on the Euclidean distance, with σ controlling the width of the kernel.

Additionally, the Gaussian kernel can be generalized by replacing the dot product $\mathbf{x}^T \mathbf{x}'$ with a nonlinear kernel function $\kappa(\mathbf{x}, \mathbf{x}')$. This leads to the extended form of the Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}')}{2\sigma^2}}$$

This extension allows the Gaussian kernel to operate in a more flexible feature space, potentially capturing nonlinear relationships between data points, thereby enhancing its applicability in various machine learning tasks.

Symbolic data kernel Kernel methods are not limited to real vectors as inputs; they can be extended to various data structures such as graphs, sets, strings, texts, and more. The kernel function serves as a measure of similarity between two samples. For example, in the case of sets, a common kernel function is employed:

$$k(A_1, A_2) = 2^{|A_1 \cap A_2|}$$

This kernel function quantifies the similarity between two sets A_1 and A_2 by computing the cardinality of their intersection. The resulting value reflects the degree of overlap between the sets, indicating their similarity.

Generative model kernel Kernel functions can also be defined using probability distributions. In the context of generative models, where $P(\mathbf{x})$ represents the probability distribution, a kernel function can be defined as:

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

This kernel function is valid as it corresponds to the inner product in a one-dimensional feature space obtained by mapping \mathbf{x} to $P(\mathbf{x})$. It effectively measures the similarity between two samples by considering their respective probabilities under the generative model.

2.4.2 Kernel ridge regression

The loss function utilized in ridge regression is given by:

$$L(\mathbf{w}) = \frac{1}{2}(\mathbf{t} - \Phi\mathbf{w})^T(\mathbf{t} - \Phi\mathbf{w}) + \frac{\lambda}{2}\mathbf{w}^T\mathbf{w}$$

To solve it, we equate the gradient of L with respect to \mathbf{w} to zero:

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \lambda\mathbf{w} - \Phi^T(\mathbf{t} - \Phi\mathbf{w}) = 0$$

Now, instead of solving it for \mathbf{w} , let's perform a variable change ($\mathbf{a} = \lambda^{-1}(\mathbf{t} - \Phi\mathbf{w})$):

$$\mathbf{w} = \Phi^T\lambda^{-1}(\mathbf{t} - \Phi\mathbf{w}) = \Phi^T\mathbf{a}$$

Substituting \mathbf{w} in the gradient, we have:

$$\begin{aligned}\lambda\mathbf{w} - \Phi^T(\mathbf{t} - \Phi\mathbf{w}) &= 0 \rightarrow \\ \Phi^T(\lambda\mathbf{a} - (\mathbf{t} - \Phi\Phi^T\mathbf{a})) &= 0 \rightarrow \\ \Phi\Phi^T\mathbf{a} + \lambda\mathbf{a} &= \mathbf{t} \rightarrow \\ \mathbf{a} &= (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{t}\end{aligned}$$

Here, $\mathbf{K} = \Phi\Phi^T$ is known as the Gram matrix. The Gram matrix is an $N \times N$ matrix where each element represents the inner product between the feature vectors:

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

The Gram matrix signifies the similarities between each pair of samples in the training data.

Prediction function To compute the prediction using the dual representation, we can utilize the following formula:

$$y(\mathbf{x}) = \mathbf{w}^T\phi(\mathbf{x}) = \mathbf{a}\Phi\phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{t}$$

Here, $\mathbf{k}(\mathbf{x})$ is defined such that $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$ for all $\mathbf{x}_n \in \mathcal{D}$. Accordingly, the prediction is computed as the linear combination of the target values of the samples in the training set.

Comparison The original representation:

- Involves computing the inverse of $(\Phi\Phi^T + \lambda\mathbf{I}_M)$, which yields an $M \times M$ matrix.
- Is computationally convenient when M is relatively small.

The dual representation:

- Requires computing the inverse of $(\mathbf{K} + \lambda\mathbf{I}_N)$, which results in an $N \times N$ matrix.
- Is computationally favorable when N is very large or even infinite.
- Eliminates the need to explicitly compute Φ , enabling application to diverse data types such as graphs, sets, strings, and text.
- The computation of the similarity between data samples (i.e., the kernel function) is typically more efficient and simpler than calculating Φ .

2.4.3 Kernel regression

The k -nearest neighbors algorithm can be utilized for regression tasks by computing the average of the target values of the k nearest samples in the training data. This can be expressed as:

$$\hat{f}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} t_i$$

Nadaraya-Watson model In k-NN regression, the model output often exhibits significant noise due to the discontinuity of neighborhood averages. The Nadaraya-Watson model, also known as kernel regression, addresses this issue by employing a kernel function to calculate a weighted average of samples:

$$\hat{f}(\mathbf{x}) = \frac{\sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) t_i}{\sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i)}$$

Typically, kernels are chosen based on their properties. Two common choices for kernels are:

- Epanechnikov Kernel (bounded support):

$$k(u) = \frac{3}{4}(1 - u^2) \quad |u| \leq 1$$

- Gaussian Kernel (infinite support):

$$K(u) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2\sigma^2}}$$

2.4.4 Gaussian processes

Starting from the assumptions of Bayesian linear regression:

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

with the following prior probability:

$$P(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \tau^2 \mathbf{I})$$

Now, let's compute the prior distribution of the outputs of the regression function:

$$\mathbf{y} = \boldsymbol{\Phi} \mathbf{w} \implies \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}, \mathbf{S})$$

Here:

- $\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}] = \boldsymbol{\Phi} \mathbb{E}[\mathbf{w}] = \mathbf{0}$
- $\mathbf{S} = \text{Cov}(\mathbf{y} \mathbf{y}^T) = \boldsymbol{\Phi} \mathbb{E}[\mathbf{w} \mathbf{w}^T] \boldsymbol{\Phi}^T = \tau^2 \boldsymbol{\Phi} \boldsymbol{\Phi}^T = \mathbf{K}$

In general, a Gaussian Process is defined as a probability distribution over a function $y(\mathbf{x})$ such that the set of values $y(\mathbf{x}_i)$ — for an arbitrary \mathbf{x}_i — jointly have a Gaussian distribution. In our case:

$$P(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})$$

where \mathbf{K} is the Gram matrix defined as:

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \tau^2 \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m)$$

This provides a probabilistic interpretation of the Kernel function as:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \mathbb{E} [y(\mathbf{x}_n), y(\mathbf{x}_m)]$$

We can apply the usual approaches to design the kernels. Two families of kernels typically used with Gaussian processes are:

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2\sigma^2}}$$

- Exponential kernel:

$$k(\mathbf{x}, \mathbf{x}') = e^{-\theta \|\mathbf{x}-\mathbf{x}'\|}$$

2.5 Support Vector Machines

Kernel methods face a notable limitation: the need to compute the kernel function for every sample in the training set. Unfortunately, this computation can be computationally infeasible in practice. To address this challenge, sparse kernel methods seek solutions that rely only on a subset of the training samples. Two well-known sparse kernel methods are:

1. Support Vector Machines (SVMs).
2. Relevance Vector Machines.

2.5.1 Separable problems

The separation between data points can also be achieved using the perceptron algorithm. However, in this case, the final result is highly dependent on the initialization.

When choosing the best solution, consider the line that separates the points. Opt for the solution with fewer points close to that separating line. To address this, we can utilize the maximum margin classifier, which computes the margin as follows:

$$\text{margin} = \min_n \frac{t_n (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

The goal is to find the optimal hyperplane by maximizing the expression:

$$\operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_n \left[\frac{t_n (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right] \right\}$$

However, solving this optimization problem can be very complex due to its computational demands and potential non-convexity.

To simplify the optimization problem, we first establish a canonical hyperplane across the separating variables. It's essential to acknowledge the existence of an infinite set of equivalent solutions represented by:

$$\kappa \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + \kappa b \quad \forall \kappa > 0$$

However, we will focus solely on solutions that adhere to the condition:

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1 \quad \forall \mathbf{x}_n \in \mathcal{S}$$

Consequently, we transform the problem into an equivalent quadratic programming task aimed at minimizing:

$$\frac{1}{2} \|\mathbf{w}\|_2^2$$

subject to the constraint $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$, for all n .

Dual problem We can obtain the dual problem by utilizing Lagrange multipliers, resulting in the following Lagrangian:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (t_i (\mathbf{w}^T \phi(\mathbf{x}_i)) - 1)$$

To maximize \mathcal{L} with respect to $\boldsymbol{\alpha}$ and minimize it with respect to \mathbf{w} and b , we compute the gradients with respect to \mathbf{w} and b and derive the dual representation:

$$\begin{cases} \frac{\partial}{\partial \mathbf{w}} \mathcal{L} = 0 \\ \frac{\partial}{\partial b} \mathcal{L} = 0 \end{cases} \rightarrow \begin{cases} \mathbf{w} = \sum_{i=1}^n \alpha_i t_i \phi(\mathbf{x}_i) \\ \sum_{i=1}^n \alpha_i t_i = 0 \end{cases}$$

This allows us to reformulate the optimization problem as the maximization of:

$$\tilde{\mathcal{L}}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to the constraints:

$$\begin{cases} \alpha_n \geq 0 \\ \sum_{n=1}^N \alpha_n t_n = 0 \end{cases} \quad \forall n = 1, \dots, N$$

where the explicit feature mapping no longer appears explicitly.

Discriminant function The resulting discriminant function can be expressed as:

$$y(\mathbf{x}) = \sum_{n=1}^N \alpha_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

Here, only samples on the margin contribute, indicated by $\alpha_i > 0$. These crucial samples are known as the Support Vectors. The bias term, denoted as b , is computed as:

$$b = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x}_n \in \mathcal{S}} \left(t_n - \sum_{\mathbf{x}_m \in \mathcal{S}} \alpha_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

This formulation ensures that the decision boundary is determined by the support vectors, reflecting the critical points in the data that define the separation between classes.

2.5.2 Non-separable problems

In our prior discussions, we've proceeded on the premise that samples are linearly separable within the feature space. Yet, this isn't universally applicable, especially in scenarios with noisy data or other complexities. To address these challenges, we introduce the concept of error (represented by ξ_i) into our classification methodology.

With this definition, we can introduce the soft-margin optimization problem, which aims to minimize:

$$\frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N \xi_n$$

subject to the constraints:

$$t_n (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1 - \xi_n \quad \forall n$$

where $\xi_n \geq 0$ are slack variables representing penalties for margin violations. The parameter C serves as a tradeoff between error and margin: it allows adjustment of the bias-variance tradeoff, and tuning may be necessary to find the optimal value for C .

Dual problem By obtaining the dual problem, we aim to maximize:

$$\tilde{\mathcal{L}}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to the constraints:

$$\begin{cases} 0 \leq \alpha_n \leq C \\ \sum_{n=1}^N \alpha_n t_n = 0 \end{cases} \quad n = 1, \dots, N$$

As usual, support vectors are the samples for which $\alpha_n > 0$. If $\alpha_n < C$, then $\xi_n = 0$, indicating that the sample is on the margin. If $\alpha_n = C$, the sample can be within the margin and either correctly classified ($\xi_n \leq 1$) or misclassified ($\xi_n > 1$).

Alternative formulation The same problem can be also formulated as the maximization of:

$$\tilde{\mathcal{L}}(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\begin{cases} 0 \leq \alpha_n \leq \frac{1}{N} \\ \sum_{n=1}^N \alpha_n t_n = 0 \\ \sum_{n=1}^N \alpha_n \geq \nu \end{cases} \quad n = 1, \dots, N$$

Where $0 \leq \nu < 1$ is a user-defined parameter that enables control over both the margin errors and the number of support vectors, ensuring that the fraction of margin errors is less than or equal to ν and the fraction of support vectors is also less than or equal to ν .

2.5.3 Support vector machines training

To solve the optimization problem and find α_i and b , several methods exist. However, the direct solution is computationally expensive, typically $O(n^3)$ where n is the size of the training set.

To mitigate this computational burden, faster approaches have been developed, including:

1. **Chunking:** Breaking the problem into smaller chunks to solve separately.
2. **Osuna's methods:** Variants of chunking methods specifically tailored for SVM optimization.
3. **Sequential Minimal Optimization (SMO):** A method that optimizes the dual problem iteratively by selecting pairs of variables to update.

Additionally, for scenarios where online learning is preferred, methods such as chunking-based approaches and incremental methods can be employed. These methods update the model gradually as new data becomes available, thus avoiding the need to retrain the entire model from scratch.

Chunking Chunking solves iteratively by addressing a sub-problem known as the working set. The working set is constructed using the current support vectors and the M samples with the largest errors (known as the worst set). It's important to note that the size of the working set may dynamically increase during the iterations. Despite this, chunking converges to the optimal solution.

Osuna's Method Osuna's method also solves iteratively by focusing on a sub-problem, the working set. However, unlike chunking, it maintains a fixed size for the working set. This method replaces some samples in the working set with misclassified samples from the dataset. Despite its fixed-size working set, Osuna's method still converges to the optimal solution.

Sequential Minimal Optimization SMO operates iteratively, but uniquely, it only works on two samples at a time. By doing so, it keeps the size of the working set minimal. Moreover, the multipliers are found analytically during each iteration. Like the other methods, SMO converges to the optimal solution.

2.5.4 Multi-class Support vector machines

One against all In the one against all approach, a k -class problem is decomposed into k binary (2-class) problems. Training involves k SVM classifiers on the entire dataset. During testing, the class selected with the highest margin among the k SVM classifiers is chosen.

One against one In one against one, a k -class problem is decomposed into $\frac{k(k-1)}{2}$ binary problems. Here, $\frac{k(k-1)}{2}$ SVM classifiers are trained on subsets of the dataset. During testing, all $\frac{k(k-1)}{2}$ classifiers are applied to the new sample, and the most voted label is chosen.

DAGSVM DAGSVM also decomposes the k -class problem into $\frac{k(k-1)}{2}$ binary problems like one-against-one. However, it employs a Direct Acyclic Graph during testing to reduce the number of SVM classifiers to apply. This leads to only $k-1$ binary SVM classifiers being involved in the test process instead of $\frac{k(k-1)}{2}$ as in one-against-one.

Summary The methods are:

- One-against-all: requires less memory but has expensive training and cheap testing.

- One-against-one: requires more memory but has slightly cheaper training and expensive testing.
- DAGSVM: moderately expensive in terms of memory requirements, with slightly cheaper training and testing.

One-against-one is considered the best performing approach due to its effective decomposition. DAGSVM provides a faster approximation of one-against-one.

Model evaluation

3.1 Bias-variance framework

The bias-variance framework provides a structured approach for evaluating model performance.

Definition (*Data*). Data are described as:

$$t_i = f(\mathbf{x}_i) + \varepsilon$$

where $\mathbb{E}[\varepsilon] = 0$ and $\text{Var}[\varepsilon] = \sigma^2$.

Definition (*Model*). The model is represented as:

$$\hat{t}_i = y(\mathbf{x}_i) + \varepsilon$$

learned from a sampled dataset $\mathcal{D} = \{\mathbf{x}_i, t_i\}$.

Definition (*Performance*). Performance is quantified by:

$$\mathbb{E}[(t - y(\mathbf{x}))^2]$$

which measures the expected squared error.

Hence, the expected squared error can be decomposed as follows:

$$\begin{aligned} \mathbb{E}[(t - y(\mathbf{x}))^2] &= \mathbb{E}[(t^2 + y(\mathbf{x})^2 - 2ty(\mathbf{x}))] \\ &= \mathbb{E}[t^2] + \mathbb{E}[y(\mathbf{x})^2] - \mathbb{E}[2ty(\mathbf{x})] \\ &= \mathbb{E}[t^2] + \mathbb{E}[t]^2 - \mathbb{E}[t]^2 + \mathbb{E}[y(\mathbf{x})^2] + \mathbb{E}[y(\mathbf{x})]^2 - \mathbb{E}[y(\mathbf{x})]^2 - 2f(\mathbf{x})\mathbb{E}[y(\mathbf{x})] \\ &= \text{Var}[t] + \mathbb{E}[t]^2 + \text{Var}[y(\mathbf{x})] + \mathbb{E}[y(\mathbf{x})]^2 - 2f(\mathbf{x})\mathbb{E}[y(\mathbf{x})] \\ &= \text{Var}[t] + \text{Var}[y(\mathbf{x})] + (f(\mathbf{x}) - \mathbb{E}[y(\mathbf{x})])^2 \\ &= \underbrace{\text{Var}[t]}_{\sigma^2} + \underbrace{\text{Var}[y(\mathbf{x})]}_{\text{variance}} + \underbrace{\mathbb{E}[f(\mathbf{x}) - y(\mathbf{x})]^2}_{\text{squared bias}} \end{aligned}$$

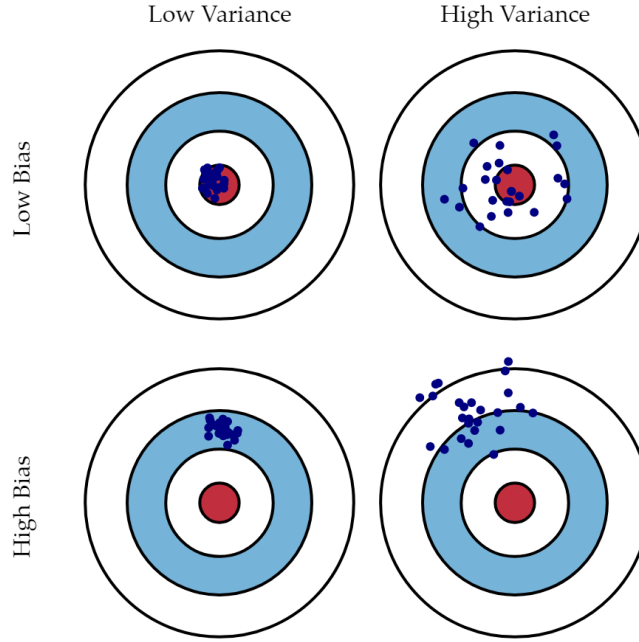


Figure 3.1: Bias-variance framework

Model variance When we sample multiple datasets \mathcal{D} , we obtain distinct models $y(\mathbf{x})$. Variance quantifies the dissimilarity between each model learned from a specific dataset and our anticipated learning outcome:

$$\text{variance} = \int \mathbb{E} [(y(\mathbf{x}) - \bar{y}(\mathbf{x}))^2] P(\mathbf{x}) d\mathbf{x}$$

The variance diminishes by simplifying the model or increasing the sample size.

Model bias Bias gauges the disparity between the truth (f) and our expected learning outcome ($\mathbb{E}[y(\mathbf{x})]$):

$$\text{bias}^2 = \int (f(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 P(\mathbf{x}) d\mathbf{x}$$

Bias decreases with more complex models.

Definition (Data noise). Data noise (σ^2) represents the variance of data and remains constant regardless of data sampling or model complexity.

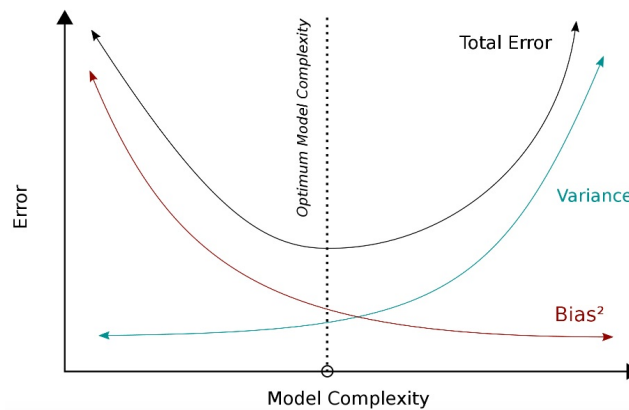


Figure 3.2: Bias-variance framework impact

In practical terms, the estimation is affected as follows:

- High variance leads to overfitting.
- High bias results in underfitting.
- Low bias and low variance yield a well-balanced approximation.

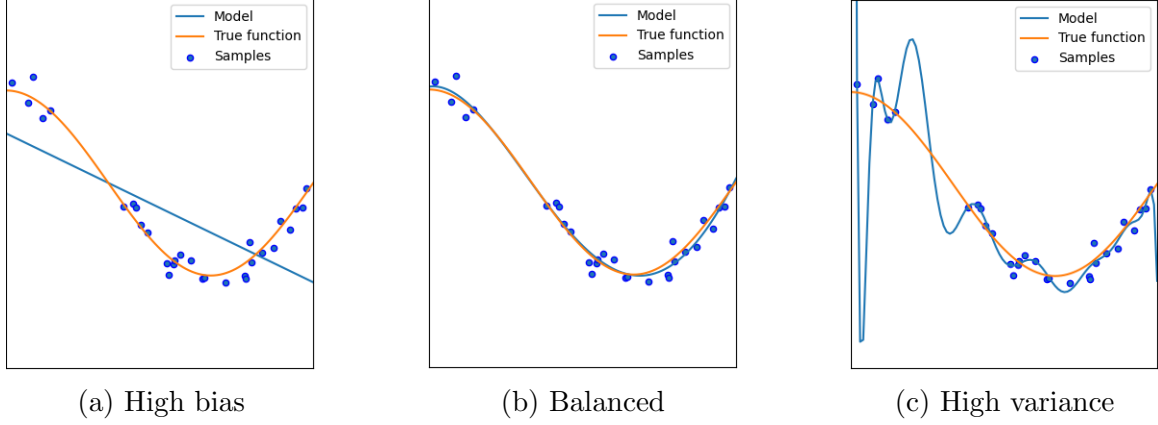


Figure 3.3: Bias-variance balancing

3.1.1 Regularization and bias-variance

The bias-variance decomposition elucidates why regularization enhances error reduction on unseen data. Lasso surpasses ridge regression when only a few features are linked to the output.

3.2 Model assessment

Given a dataset $\mathcal{D} = \{\mathbf{x}_i, t_i\}$ with $i = 1, \dots, N$, we can choose a model based on the computed loss L on \mathcal{D} . For regression, the loss function is defined as:

$$L_{train} = \frac{1}{N} \sum_{n=1}^N (t_n - y(\mathbf{x}_n))^2$$

And for classification, the loss function becomes:

$$L_{train} = \frac{1}{N} \sum_{n=1}^N I(t_n \neq y(\mathbf{x}_n))$$

The training error decreases as the model complexity increases.

However, it's important to note that the training error doesn't give an accurate estimate of the error on new data, known as the prediction error. For regression, the prediction error is represented as:

$$L_{true} = \iint (t - y(\mathbf{x}))^2 P(\mathbf{x}, t) d\mathbf{x} dt$$

And for classification, it is:

$$L_{true} = \iint I(t \neq y(\mathbf{x}))P(\mathbf{x}, t)d\mathbf{x}dt$$

Unfortunately, modeling the joint probability distribution $P(\mathbf{x}, t)$ is often not feasible.

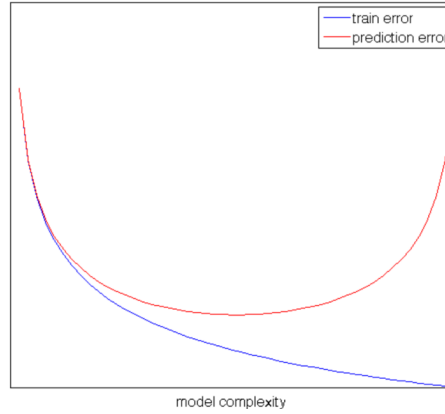


Figure 3.4: Train error compared to prediction error

Practical application In practical scenarios, data is typically randomly split into a training set and a test set. Model parameters are optimized using the training set, and the prediction error is estimated using the test set. For regression, this estimation yields:

$$L_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (t_n - y(\mathbf{x}_n))^2$$

And for classification:

$$L_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} I(t_n \neq y(\mathbf{x}_n))$$

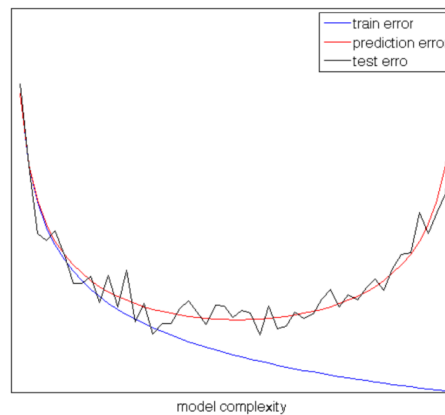


Figure 3.5: Error in practice

As the number of data points increases, these errors tend to converge, as depicted in the following figure:

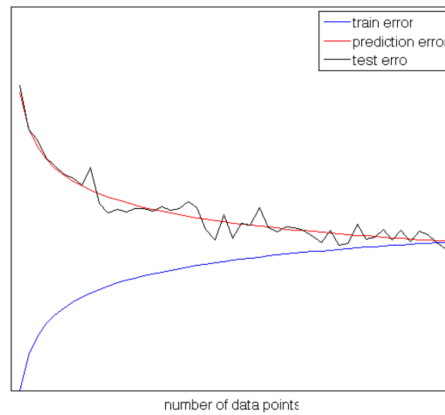


Figure 3.6: Error in function of the number of data points

Analyzing the train-test errors helps identify potential issues:

- *High bias*: when both training and test errors are higher than expected and close to each other.
- *High variance*: when the training error is significantly lower than expected and gradually approaches the test error.

Problems Frequently, data availability is constrained, and the test error tends to be minimal, leading to potential overestimation or underestimation of prediction error. Utilizing test error for model selection can result in overfitting to the test set. An unbiased estimate of prediction error is achievable only if the test set remains separate from both training and model selection phases.

3.2.1 Optimal model

To select the optimal model and determine the appropriate hyperparameters, we initially partition the data into three subsets: training data, validation data, and test data. The process involves the following steps:

1. Utilize the training data to train the model parameters.
2. For each trained model, assess its performance using the validation data to compute the validation error.
3. Identify the model with the lowest validation error, and subsequently, employ the test data to estimate the prediction error.

However, for this approach to be dependable, it's imperative that the validation data set is sufficiently sizable, especially in comparison to the training data set. Otherwise, there's a risk of overfitting to the validation data, potentially leading to the selection of a suboptimal model.

Leave-one-out cross validation Leave-one-out cross-validation (LOO-CV) involves training the model on all samples in the dataset \mathcal{D} except for a single sample $\{\mathbf{x}_i, t_i\}$, and then

evaluating the model's performance on that omitted sample. The prediction error estimate of our model is then computed as the average error across all single-sample evaluations:

$$L_{LOO} = \frac{1}{N} \sum_{i=1}^N (t_i - y_{\mathcal{D}_i}(\mathbf{x}_i))^2$$

Here, $y_{\mathcal{D}_i}$ represents the model trained on \mathcal{D} excluding $\{\mathbf{x}_i, t_i\}$.

The L_{LOO} estimate of prediction error provides an almost unbiased assessment (slightly pessimistic). However, LOO-CV is computationally intensive due to its requirement to repeatedly train models on nearly all data points.

K-fold cross validation K-fold cross-validation involves randomly dividing the training data \mathcal{D} into k folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$. For each fold \mathcal{D}_i , the model is trained on \mathcal{D} excluding \mathcal{D}_i , and then the error is computed on \mathcal{D}_i as follows:

$$L_{\mathcal{D}_i} = \frac{k}{N} \sum_{(\mathbf{x}_n, t_n) \in \mathcal{D}_i} (t_n - y_{\mathcal{D} \setminus \{\mathcal{D}_i\}}(\mathbf{x}_n))^2$$

Finally, the prediction error is estimated as the average error computed across all folds:

$$L_{k-fold} = \frac{1}{k} \sum_{i=1}^k L_{\mathcal{D}_i}$$

The L_{k-fold} estimate of prediction error provides a slightly biased (pessimistic) assessment but is computationally less expensive. Typically, k is set to ten.

Other metrics Various metrics are employed to evaluate models by adjusting their training error based on their complexity:

- Mallows's C_p :

$$C_p = \frac{1}{N} (\text{RSS} + 2M\sigma^2)$$

- Akaike Information Criteria:

$$\text{AIC} = -2 \ln(L) + 2M$$

- Bayesian Information Criteria:

$$\text{BIC} = -2 \ln(L) + M \ln(N)$$

- Adjusted R^2 :

$$A_{R^2} = 1 - \frac{\text{RSS}/(n - m - 1)}{\text{TSS}/(N - 1)}$$

Here, M represents the number of parameters, N denotes the number of samples, L signifies the loss function, σ^2 stands for the estimate of noise variance, RSS corresponds to the residual sum of squares, and TSS indicates the total sum of squares.

AIC and BIC are typically utilized when maximizing the log-likelihood. BIC tends to penalize model complexity more severely compared to AIC.

3.3 Model complexity

Introducing an additional feature leads to an exponential growth in the volume of the input space. This growth leads to the following problems:

- *Computational cost*: the computational resources required to process and analyze the expanded input space increase significantly.
- *Data quantity*: the amount of data needed to effectively explore and train models in the expanded input space may be substantial.
- *Large model variance* (overfitting): with the increased complexity of the input space, there is a higher risk of models capturing noise or irrelevant patterns, leading to overfitting and decreased generalization performance.

Our goal is to choose the model with the minimal prediction error, which can be attained by decreasing the variance of the model:

- *Feature selection*: by carefully designing the feature space, we can choose the most impactful subset from all available features.
- *Dimensionality reduction*: mapping the input space to a lower-dimensional representation can effectively reduce complexity and variance.
- *Regularization*: shrinkage of parameter values towards zero helps control model complexity and mitigate overfitting.

These approaches are not mutually exclusive and can be combined to enhance model performance.

3.3.1 Feature selection

The most straightforward approach appears to be comparing all possible combinations of features. Given M features, for each $k = 1, \dots, M$, we would need to train all $\binom{M}{k} = \frac{M!}{k!(M-k)!}$ models with exactly k features and select the optimal one. However, this procedure quickly becomes computationally impractical.

In practical scenarios, feature selection is often carried out based on the specific model being utilized:

- *Filter*: features are assessed individually using certain evaluation metrics (e.g., correlation, variance, information gain), and the top k features are selected. While this method is very fast, it fails to capture any subset of mutually dependent features.
- *Embedded*: feature selection is integrated into the machine learning approach itself (e.g., lasso, decision trees). Although this method is not computationally expensive, the features identified are specific to the chosen learning approach.
- *Wrapper*: a search algorithm is employed to identify a subset of features by iteratively training a model with different feature subsets and evaluating their performance. This method utilizes either a simpler model or a basic machine learning approach to evaluate the features. Greedy algorithms are typically employed to search for the best feature subset.

3.3.2 Dimensionality reduction

Dimensionality reduction aims to decrease the dimensions of the input space, but it differs from feature selection in two significant ways:

- It utilizes all features and transforms them into a lower-dimensional space.
- It is an unsupervised approach, meaning it doesn't rely on labeled data for training.

There are numerous methods for performing dimensionality reduction, including:

- Principal Component Analysis (PCA).
- Independent Component Analysis (ICA).
- Self-Organizing Maps.
- Autoencoders.
- ISOMAP.
- t-SNE.

3.4 Ensemble

We've explored methods to decrease variance while balancing increased bias. However, we want to reduce variance without amplifying bias or mitigate bias altogether.

These objectives can indeed be achieved through the utilization of two ensemble methods involving the learning of multiple models and their combination:

- *Bagging*: involves training multiple models independently on different subsets of the data and then combining their predictions.
- *Boosting*: utilizes an iterative approach where models are sequentially trained, each aiming to correct the errors of its predecessors, leading to the creation of a strong ensemble model.

3.4.1 Bagging

Let assume to have N datasets and to learn from them N models, y_1, y_2, \dots, y_N . Now let us compute an aggregate model as:

$$y_{AGG} = \frac{1}{N} \sum_{i=1}^N y_i$$

If the datasets are independent, the model variance of y_{AGG} will be $\frac{1}{N}$ of the model variance of the single model y_i . However, we generally do not have N datasets.

Bagging, short for Bootstrap Aggregation, involves the following steps:

1. Generate N datasets by applying random sampling with replacement.
2. Train a model (classification or regression) using each dataset generated.

3. To predict new samples, apply all the trained models and combine their outputs using majority voting (for classification) or averaging (for regression).

Bagging is generally beneficial as it reduces variance, although the sampled datasets are not independent. It proves particularly useful with unstable learners, characterized by significant changes with small dataset variations (low bias and high variance), and in scenarios with a high degree of overfitting (low bias and high variance). However, it does not offer much help with robust learners, which are insensitive to data changes (typically higher bias but lower variance).

3.4.2 Boosting

Boosting aims to minimize bias by employing a series of simple (weak) learners.

The core concept of boosting involves iteratively training a sequence of weak learners, with each iteration concentrating on the samples misclassified in the preceding iteration.

Ultimately, an ensemble model is constructed by combining the outputs of all the weak learners trained.

3.4.3 Summary

The characteristics of bagging are:

- Decreases variance.
- Less effective for stable learners.
- Applicable with noisy data.
- Generally beneficial, though the improvement might be modest.
- Naturally suited for parallelization.

The characteristics of boosting are:

- Reduces bias (typically without overfitting).
- Compatible with stable learners.
- May encounter challenges with noisy data.
- Not always effective, but can yield significant improvements.
- Sequential in nature.

APPENDIX A

Algebra and statistics

A.1 Least squares

Let's reconsider a dataset consisting of N inputs $\mathbf{x}_i = (x_{i1}, \dots, x_{iD})$, where each $x_{ij} \in \mathbb{R}$ represents a feature with dimension D , along with a target $t_i \in \mathbb{R}$ associated with each input \mathbf{x}_i .

Our aim is to predict the target t_i by computing a linear combination of the input \mathbf{x}_i , which involves generating a parameter vector $\mathbf{w} = (w_1, \dots, w_D)^T$ that minimizes a certain loss function. Specifically, if we consider the loss function to be the summation of squared prediction errors, it corresponds to Least Square minimization.

Now, let's define the following loss function:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(t_i + \sum_{j=1}^D x_{ij} w_j \right)^2$$

By defining the matrix $X = [\mathbf{x}_1 \ \dots \ \mathbf{x}_N]^T$, we can rewrite the loss function as:

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - X\mathbf{w}\|_2^2 = \frac{1}{2} (\mathbf{t} - X\mathbf{w})^T (\mathbf{t} - X\mathbf{w})$$

Minimizing the Loss To minimize the loss, we need to compute its derivatives with respect to each component of \mathbf{w} :

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \left(\frac{\partial L(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial L(\mathbf{w})}{\partial w_D} \right)$$

In this case, we have two different formulations for the derivative:

1. Traditional (element-wise) derivative:

$$\left(\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} \right)_h = \frac{\partial L(\mathbf{w})}{\partial w_h} = \frac{\partial}{\partial w_h} \left[\frac{1}{2} \sum_{i=1}^N \left(t_i - \sum_{j=1}^D x_{ij} w_j \right)^2 \right]$$

2. Matrix derivative:

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \left[\frac{1}{2} (\mathbf{t} - X\mathbf{w})^T (\mathbf{t} - X\mathbf{w}) \right]$$

A.2 Matrices

Eigenvalues and eigenvectors For a square matrix $\mathbb{R}^{n \times n}$, the corresponding eigenvector equations are given by:

$$A\mathbf{v}_i = \lambda_i \mathbf{v}_i$$

Here:

- The eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ represent directions unaffected by the transformation A .
- The eigenvalues $\lambda_1, \dots, \lambda_n$ determine the scaling factor for the corresponding eigenvectors \mathbf{v}_i .

In matrix notation, we can express this relationship as:

$$(A - \lambda I_n)\mathbf{v} = 0$$

This equation has a non-trivial solution only if the rank of the matrix $A - \lambda I_n$ is full, or equivalently:

$$|A - \lambda I_n| = 0$$

Property A.2.1. The rank of A is equal to the number of non-zero eigenvalues.

Property A.2.2. The determinant of A is equal to the product of its eigenvalues:

$$|A| = \prod_{i=1}^n \lambda_i$$

Trace The trace of A , denoted as $\text{Tr}(A)$, is equal to the sum of its eigenvalues:

$$\text{Tr}(A) = \sum_{i=1}^n \lambda_i$$

A.2.1 Properties

Definition (*Positive definite matrix*). A matrix A is said to be positive definite if $\mathbf{x}^T A \mathbf{x} > 0$ for all vectors $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$.

A positive definite matrix has all positive eigenvalues, i.e., $\lambda_i > 0$ for all i .

Definition (*Semi-positive definite matrix*). A matrix A is said to be semi-positive definite if $\mathbf{x}^T A \mathbf{x} \geq 0$ for all vectors $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$.

A semi-positive definite matrix has all non-negative eigenvalues, i.e., $\lambda_i \geq 0$ for all i .

A.3 Random variables

A discrete random variable X is a variable with values in a discrete set, whose value is determined by a stochastic phenomenon. We define a probability function $P : E \rightarrow [0, 1]$ which indicates the likelihood of events in E :

$$P(X = i) = \frac{|i|}{|E|}$$

A properly defined probability function should satisfy the following properties:

1. $\forall i \in E, 0 \leq P(X = i) \leq 1.$
2. $\sum_{i \in E} P(X = i) = 1.$

Cumulative function Assuming events are ordered, a cumulative function $F : E \rightarrow [0, 1]$ defines the probability of multiple events:

$$F(i) = P(X \leq i) = \sum_{h=1}^i P(X = h) = \sum_{h \in E, h \leq i} \frac{|h|}{|E|}$$

A properly defined cumulative function should satisfy the following properties:

- $0 \leq F(i) \leq 1$
- $F(i) = 0$ for all $i < \min_{h \in E} h$
- $F(i) = 1$ for all $i \geq \max_{h \in E} h$

Random variables characteristics Quantities characterizing a random variable include the expected value (moment of order one):

$$\mathbb{E}[X] = \sum_{i \in E} iP(X = i)$$

And the variance (moment of order two):

$$\text{Var}(X) = \sum_{i \in E} (\mathbb{E}[X] - i)^2 P(X = i)$$

The standard deviation of the variance is defined as:

$$\text{std}(X) = \sqrt{\text{Var}(X)}$$

A.3.1 Continuous random variable

Similarly, if the set E is not discrete, we could define the probability density function as:

$$f(x) = \lim_{\delta x \rightarrow 0} \frac{P(x \leq X \leq x + \delta x)}{\delta x}$$

The properties of continuous random variables are:

- $f(x) \geq 0$ for all $x \in \Omega$
- $\int_{x \in \Omega} f(x) dx = 1$

Cumulative density function The cumulative density function is defined as:

$$F(x) = \int_{s \in \Omega, s \leq x} f(s) ds$$

The cumulative density function has the following properties:

- $0 \leq F(x) \leq 1$ for all $x \in \Omega$
- $F(\min_{x \in \Omega} x - \varepsilon) = 0$ for all $\varepsilon > 0$
- $F(\max x \in \Omega) = 1$

A quantile of order α is defined as a point $z_\alpha \in \Omega$ such that:

$$F(z_\alpha) = 1 - \alpha$$

or a point of the domain leaving to its left a cumulated probability of $1 - \alpha$.

Random variables characteristics Quantities which characterize a random variable are the expected value (moment of order one):

$$\mu = \mathbb{E}[X] = \int_{x \in \Omega} x f(x) dx$$

And the variance (moment of order two):

$$\sigma^2 = \text{Var}(X) = \int_{x \in \Omega} (\mathbb{E}[X] - x)^2 f(x) dx$$

Gaussian distribution The Gaussian distribution is expressed as:

$$f(x, \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

And the relative cumulative density function:

$$F(x, \mu, \sigma) = \int_{-\infty}^x f(t, \mu, \sigma) dt$$

A.4 Distributions

Usually, we do not have any information about the distribution of random variables. Therefore, we need to estimate their mean and variance. A consistent estimator for the expected value is:

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$$

A consistent estimator for the variance is:

$$\bar{s} = \frac{\sum_{i=1}^n (\bar{X} - x_i)^2}{n - 1}$$

Theorem A.4.1 (Central limit). *Assuming $\{X_1, \dots, X_n\}$ as a sequence of independent and identically distributed random variables with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2 < \infty$, then:*

$$\sqrt{n} \left(\frac{\sum_{i=1}^n X_i}{n} - \mu \right) \rightarrow \mathcal{N}(0, \sigma^2)$$

where the convergence holds in distribution.

A.5 Confidence intervals

We need to establish a level of confidence to determine if our estimator is sufficiently accurate. Since the probability of $\bar{X} = \mathbb{E}[X]$ is zero, given that the realization of the expected value is a continuous random variable itself, we need to construct intervals where we have high confidence that the true mean $\mathbb{E}[X]$ lies within. A 95% confidence interval implies that it works correctly 95% of the time.

The available options for confidence intervals are:

- Gaussian approximation:

$$\bar{X} - \frac{z_{\frac{\alpha}{2}}\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{z_{\frac{\alpha}{2}}\sigma}{\sqrt{n}}$$

- Chebyshev's inequality (requires $\mathbb{E}[X] = \mu < \infty$ and $\text{Var}[X] = \sigma^2 < \infty$):

$$\bar{X} - \frac{\sigma}{\sqrt{n}\sqrt{\alpha}} \leq \mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}\sqrt{\alpha}}$$

- Hoeffding's inequality (finite support):

$$\bar{X} - (b-a)\sqrt{\frac{-\log(\frac{\alpha}{2})}{2n}} \leq \mu \leq \bar{X} + (b-a)\sqrt{\frac{-\log(\frac{\alpha}{2})}{2n}}$$

A.6 hypothesis testing

In hypothesis testing, we aim to demonstrate that the estimated parameter \bar{X} is equal to μ and that another estimated parameter \bar{X}' is different from yet another estimated parameter \bar{X}'' . With statistics, we define a null hypothesis H_0 and an alternative hypothesis H_1 :

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_1 : \mu \neq \mu_0$$

and utilize the data to provide evidence supporting either hypothesis.

A.6.1 Basic Gaussian test

Given the data $\{x_1, \dots, x_n\}$, we have:

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \rightarrow t = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim \mathcal{N}(0, 1)$$

Fixing a confidence $1 - \alpha$ with $\alpha \in (0, 1)$, the test statistic t should be close to the true mean μ with very high probability. Formally:

$$\mathbb{P}(t < z_{\frac{\alpha}{2}} \vee t > z_{1-\frac{\alpha}{2}}) = \alpha$$

The corresponding decision table is:

		Decision	
		Fail to reject H_0	Reject H_0
True	H_0	Correct	Type I error (α)
	H_1	Type II error	Correct

A.6.2 P-value

To avoid specifying the confidence α , we can let the data inform us about how confident we might be about their correspondence to a specific hypothesis:

- Small p-values imply that we are confident that the H_1 hypothesis holds.
- Large p-values imply that we are not able to reject the H_0 hypothesis.

A.7 Bayesian approach

The bounds provided by traditional methods do not allow for the incorporation of information one has about the distribution parameters. A new approach considers the expected value μ of the random variable X as a random variable itself. Bayes' formula is utilized to update this information:

$$P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

Considering Bayes' formula, we have:

$$\begin{aligned} P(\mu|x_1, \dots, x_t) &= \frac{P(x_1, \dots, x_t|\mu)P(\mu)}{P(x_1, \dots, x_t)} \\ &\propto P(x_t|\mu)P(x_1, \dots, x_{t-1}|\mu)P(\mu) \\ &= P(x_t|\mu)P(x_{t-1}|\mu)P(x_1, \dots, x_{t-2}|\mu)P(\mu) \\ &= P(\mu) \prod_{h=1}^t P(x_h|\mu) \end{aligned}$$

We incrementally incorporate information from a prior distribution $P(\mu)$.

APPENDIX B

Python introduction
