

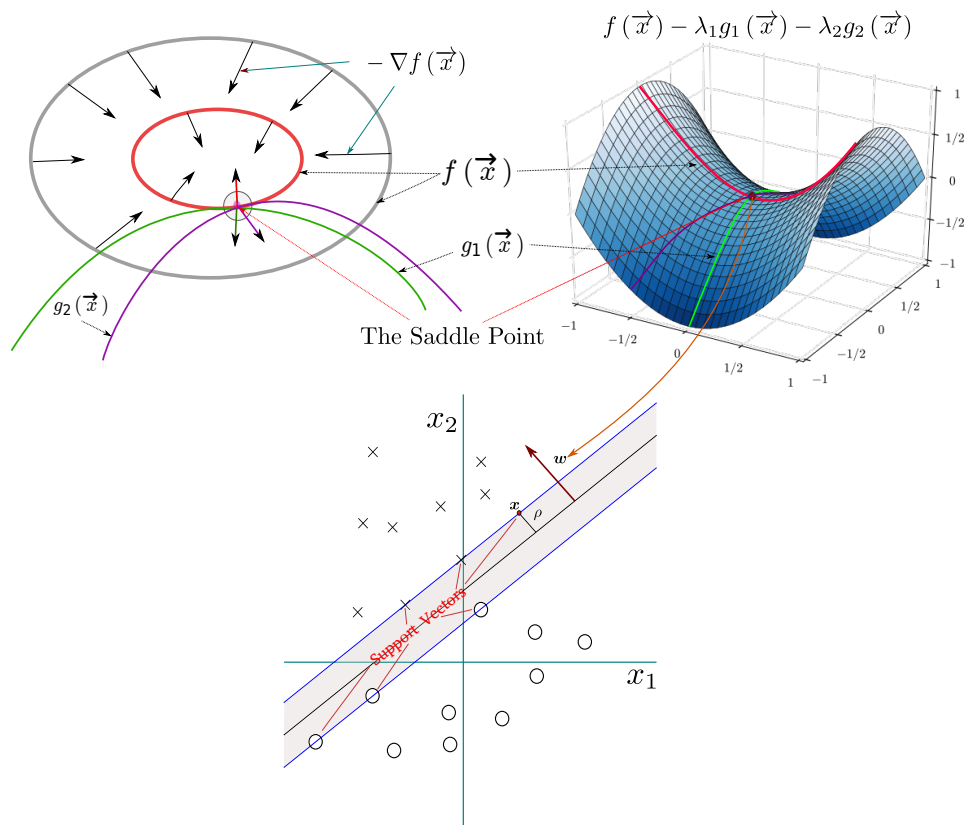
Introduction to Machine Learning

Notes

A Basic Introduction to Learning

Andres Mendez-Vazquez

June 6, 2018



Contents

1	Introduction	3
2	Loss Functions and Empirical Risk	4
2.1	Empirical Risk as Error Minimization	4
2.2	Example of Risk Minimization, Least Squared Error	6
3	Two Methods of Learning	7
3.1	Least Squared Error	8
3.2	Nearest Neighborhood	10
3.3	Machine Learning Methods as combination and improvements of LSE and k -NN	11
4	Supervised Learning as Function Approximation	12
4.1	Regression as Controlled Overfitting	12
4.2	Extreme Fitting, Bias-Variance Problem	14
4.2.1	“Extreme” Cases of Fitting	16
5	Some Classes of Estimators	18
5.1	Roughness Penalty	18
5.2	Kernel Methods and Local Regression	19
5.3	Basis Functions and Dictionary Methods	20
6	Conclusion	20

List of Figures

1	Mapping from the input space into the output space with two classes	3
2	A Hyperplane splitting the space in	4
3	The positive subspace \mathcal{R}_1 and the negative subspace \mathcal{R}_2	5
4	ϵ as the difference $\epsilon = y - g(\mathbf{x}, \mathbf{w})$	9
5	Decision Boundary $\left\{ \mathbf{x} \mid \hat{\mathbf{w}}^T \mathbf{x} = (l_1 + l_2)/2 \right\}$	10
6	Mahalanobis as a distance measuring the maximum change in each feature.	11
7	The Decision Boundary $k = 5$ and $k = 1$. Here $k=1$ correspond to a Voronoi diagram [1].	11
8	The data and the random observations.	13
9	The function used to estimate the output y is highly biased given the mapping.	16
10	The high variance case of fitting	17
11	The cubic smooth spline	18
12	Two types of Kernels	19

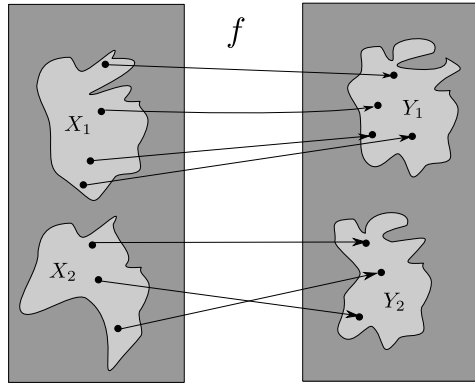


Figure 1: Mapping from the input space into the output space with two classes

1 Introduction

As in any other field of Science, it is necessary to obtain a basic understanding of the subject before getting into the main matter of the subject. In our particular case, the concept of “Learning” in Machine Learning [2, 3, 4, 5]. As we will see through our study of Machine Learning, we could give the following definition of “Learning.”[6]

Definition 1.1. *Given that the information of an object has been summarized by d features comprised as a feature vector $\mathbf{x} \in \mathbb{R}^d$, and each of these objects has been labeled by elements in a set $\{y_i \in \mathbb{R}\}$. This allows to split the set of object into a series classes, as for example $y_i \in \{-1, 1\}$. Then, for example, the process of learning is the generation of a mapping $f : \mathbb{R}^d \mapsto \{y_i\}$ (Figure 1) such that the squared error estimation of the class label of a new sample is minimized (Equation 1).*

$$\min_{\hat{f}} E_{\mathcal{X}, \mathcal{Y}} \left[\left(\hat{f}(\mathbf{x}) - y \right)^2 \mid \mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d, y \in \mathcal{Y} \subseteq \mathbb{R} \right] \quad (1)$$

Although this task looks quite simple, actually it can be quite complex! Therefore, a good starting point to this endeavors is to use our analytic geometry [7] from high school to try to obtain such mapping. Thus, if we use a geometric approach, and assuming two linearly separable classes, we could start by splitting the samples into two set of elements,

$$\begin{aligned} C_1 &= \{\mathbf{x}_i \in \mathbb{R}^d \mid y_i = 0\} \\ C_2 &= \{\mathbf{x}_i \in \mathbb{R}^d \mid y_i = 1\} \end{aligned}$$

Then, we could try to obtain a geometric function allowing to split the space of inputs into two classes. For example, we could use (Figure 2) .

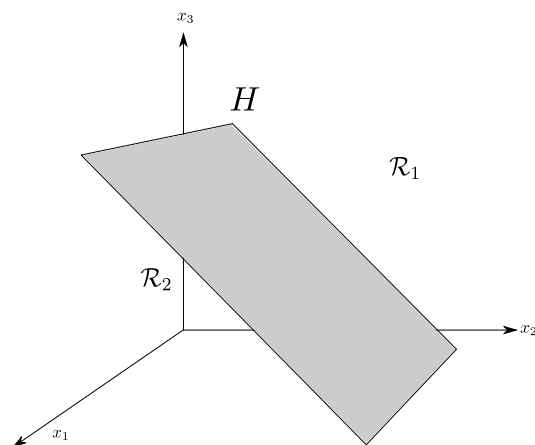


Figure 2: A Hyperplane splitting the space in

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (2)$$

Question, How does this function split the space in \mathbb{R}^3 ? The answer comes from the equation (Equation 3).

$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos \theta \quad (3)$$

Basically, θ represents the angle between vectors \mathbf{w} and \mathbf{x} . Thus, we have the following rule:

1. If $0 \leq \theta \leq 90$ then $\mathbf{w}^T \mathbf{x} \geq 0$
2. If $\theta > 90$, then $\mathbf{w}^T \mathbf{x} < 0$

Therefore, \mathbf{w} defines a positive subspace and a negative subspace (Figure 3) given the direction defined by the vector.

A question arises naturally, from this first approach, How do you obtain such function? We will see a first approach for this in the following section.

2 Loss Functions and Empirical Risk

2.1 Empirical Risk as Error Minimization

The first thing that you need to think about, while trying to obtain an estimation mapping between the input space \mathcal{X} and the output space \mathcal{Y} , is our decision about how to relate the two variables, the vectors $\mathbf{x} \in \mathbb{R}^d$ and the outputs $y \in \mathbb{R}$.

Example 2.1. We could try to relate the variations in prices at supermarkets with the variations in salary by assuming a growing economy [8]. Given such scenario, you could think a linear relation like (Equation 4).

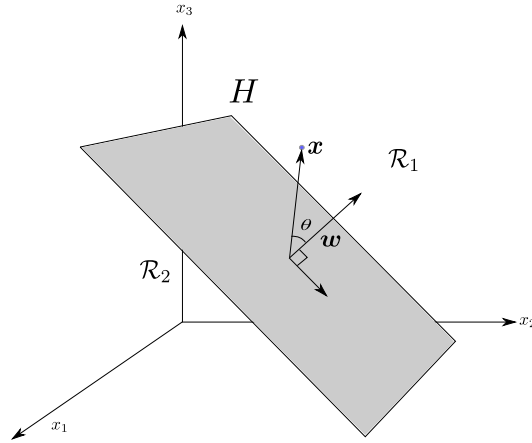


Figure 3: The positive subspace \mathcal{R}_1 and the negative subspace \mathcal{R}_2

$$\Delta P = w_0 + w_1 \Delta S + \epsilon \quad (4)$$

with ΔP and ΔS are the price and salary variations, w_0 and w_1 are unknown coefficients and $\epsilon \sim N(0, \sigma^2)$.

This example depicts one of the classic models in the literature, the linear models where a random source of noise is used to model the variability on the y_i being observed. Now, given the previous linear random relation, anybody will start to wonder how to minimize the noise or error (Equation 5),

$$\epsilon = \Delta P - (w_0 + w_1 \Delta S). \quad (5)$$

Given that we want to minimize such error, we require to compose the previous functions with an specific one to obtain the desired results. For this, we introduce the idea of statistical risk as a functional that evaluates the expected value of a function that describe a measure of loss.

Definition 2.2. (Principle of Empirical Risk) [9] Given a sequence of data samples, x_1, x_2, \dots, x_N sampled iid from a distribution $P(x|\Theta)$, and an hypothesis function $h : X \mapsto Y$ that allows to map the samples x_i into a particular output y_i . A measure of the risk of missing the estimation, $h(x)$, is found by using a function, called loss function, measuring the difference between the desired output y_i and the estimation $h(x_i)$. Thus, the Empirical Risk is defined as the expected value of the loss function based in the joint distribution $P(x, y)$ (Equation 6).

$$R(h) = E_{X,Y} [\ell(h(x), y)] = \int_{X,Y} \ell(h(x), y) p(x, y) dx dy \quad (6)$$

Examples of such loss functions are (Equation 7):

$$\begin{aligned}
 h(w_0, w_1) &= [\Delta P - (w_0 + w_1 \Delta S)]^2 \quad (\text{Least Squared Error}), \\
 h(w_0, w_1) &= |\Delta P - (w_0 + w_1 \Delta S)| \quad (\text{Absolute Difference}), \\
 h(w_0, w_1) &= \frac{1}{1 + \exp\{-(w_0 + w_1 \Delta S)\}} \quad (\text{Sigmoidal}).
 \end{aligned} \tag{7}$$

After looking at such functions loss functions, one could ask, Is there a way to obtain the optimal function h that minimizes their Empirical Risks? Is there a way to find the optimal loss function for a given data set $\{(x_i, y_i)\}$? These simple questions have driven the Machine Learning endeavor for the last 60 years.

2.2 Example of Risk Minimization, Least Squared Error

To exemplify the idea of minimizing the Risk Function [3], let us to select a convenient loss function:

$$L(y, f(x)) = (y - f(x))^2 \quad (\text{Least Squared Error}) \tag{8}$$

Here, we select a functional relation between the outputs y 's and inputs x 's:

$$Y_{noisy}(x) = f(y) + \epsilon \quad \text{with } \epsilon \sim N(0, 1) \tag{9}$$

Then, we have the following Risk functional [9, 3]:

$$R(f) = E(Y - f(X))^2 = \int_X \int_Y [y - f(x)]^2 p_{xy}(x, y) dx dy$$

Now, we can condition the probability density function with respect to X :

$$p(X, Y) = p(Y|X)p(X)$$

Therefore, we have

$$\begin{aligned}
 \int_X \int_Y [y - f(x)]^2 p_{xy}(x, y) dx dy &= \int_X \int_Y [y - f(x)]^2 p_{y|x}(y|x) p_x(x) dx dy \\
 &= \int_X \left[\int_Y [y - f(x)]^2 p_{y|x}(y|x) dy \right] dx \\
 &= E_X \left[\int_Y [y - f(x)]^2 p_{y|x}(y|x) dy \right] \\
 &= E_X E_{Y|X} \left[(Y - f(X))^2 | X \right]
 \end{aligned}$$

Thus, if we fix $X = \mathbf{x}$, we have that the Risk functional is with respect to such \mathbf{x} :

$$R(f)_{X=\mathbf{x}} = E_{Y|X=\mathbf{x}} \left[(Y - f(\mathbf{x}))^2 | X = \mathbf{x} \right]$$

This function can be optimized by realizing that

$$\begin{aligned}
 E_{Y|X=\mathbf{x}} \left[(Y - f(\mathbf{x}))^2 | X = \mathbf{x} \right] &= E_{Y|X=\mathbf{x}} \left[(Y - \bar{Y})^2 | X = \mathbf{x} \right] + \dots = \\
 &E_{Y|X=\mathbf{x}} \left[(\bar{Y} - f(\mathbf{x}))^2 | X = \mathbf{x} \right] + \dots \\
 &2(\bar{Y} - f(\mathbf{x})) E_{Y|X=\mathbf{x}} \left[(Y - \bar{Y}) | X = \mathbf{x} \right]
 \end{aligned}$$

Then, you notice that

$$\begin{aligned}
 E_{Y|X=\mathbf{x}} \left[(Y - \bar{Y}) | X = \mathbf{x} \right] &= E_{Y|X=\mathbf{x}} [Y] - E_{Y|X=\mathbf{x}} \left[\frac{1}{N} \sum_{i=1}^N Y_i \right] \\
 &= \mu_Y - \frac{1}{N} \sum_{i=1}^N E_{Y|X=\mathbf{x}} [Y_i] \\
 &= \mu_Y - \frac{N\mu_Y}{N} \\
 &= 0
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 E_{Y|X=\mathbf{x}} \left[(Y - f(\mathbf{x}))^2 | X = \mathbf{x} \right] &= E_{Y|X=\mathbf{x}} \left[(Y - \bar{Y})^2 | X = \mathbf{x} \right] + \dots \\
 &E_{Y|X=\mathbf{x}} \left[(\bar{Y} - f(\mathbf{x}))^2 | X = \mathbf{x} \right]
 \end{aligned}$$

Therefore, if we choose

$$f(X) = \bar{Y} \approx E_Y [Y|X = \mathbf{x}]$$

Thus, we finish with

$$E_{Y|X=\mathbf{x}} \left[(Y - f(\mathbf{x}))^2 | X = \mathbf{x} \right] = E_{Y|X=\mathbf{x}} \left[(Y - E_Y [Y|X = \mathbf{x}])^2 | X = \mathbf{x} \right] = \sigma_Y^2 = 1,$$

the variance of $Y \sim N(E_Y [Y|X = \mathbf{x}], 1)$ when choosing $f(X) = E_Y [Y|X = \mathbf{x}]$ given $\epsilon \sim N(0, 1)$. This means that the optimal estimator is the conditional mean for a point $X = \mathbf{x}$, when measured by the expected squared error.

3 Two Methods of Learning

Therefore, given the Empirical Risk, we can think on the process of “learning” as the process of finding an optimal function h that minimizes such risk under an specific loss function/learning model. Clearly, this process does not take in account an important issue in the human process of learning:

- How a human being uses few samples to generalize a concept?

For example, once a child learns the concept of “cat,” she/he will be able to recognize almost every cat over the planet. Although, the concept of generalization exists in Machine Learning as part of the Bias-Variance problem [2, 3, 5], little progress has been done toward a sparse inductive generalization. Nevertheless,

there have been claims in the Deep Neural Network field of new generalization powers [10] by those devices. However, those powers require massive amounts of data, unlikely as the human child learning something as not so simple as a “cat.”

After the previous digression, it is time to go back to our main topic, the process of “Learning” in Machine Learning. For this, we will take a look at two different methods of learning that could be seen as the extremes of a line where the Machine Learning methods live.

3.1 Least Squared Error

Here, we have the following situation about the data samples of a given classification problem:

- A sequence of samples with their respective labels, $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- Here, each $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ with $i = 1, 2, \dots, N$.

Then, it is possible to use the squared loss function, under a function g estimating the outputs y_i 's:

$$\ell(g(\mathbf{x}), y) = (g(\mathbf{x}) - y)^2 \quad (10)$$

Thus, we have:

$$R(g) = \int_{X,Y} \ell(g(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x}dy = \sum_{i=1}^N (g(\mathbf{x}_i) - y_i)^2 p(\mathbf{x}_i, y_i). \quad (11)$$

Now, if assume two classes at the data i.e. $y_i \in \{l_1, l_2 | l_i \in \mathbb{R}\}$ such that $|C_{l_1}| = |C_{l_2}|$

$$p(\mathbf{x}_i, y_i) = \frac{|C_{l_1}|}{N} \approx \frac{1}{2} \text{ if } \mathbf{x}_i \in C_{l_1},$$

$$p(\mathbf{x}_i, y_i) = \frac{|C_{l_2}|}{N} \approx \frac{1}{2} \text{ if } \mathbf{x}_i \in C_{l_2}.$$

Therefore,

$$R(g) = \sum_{i=1}^N (g(\mathbf{x}_i) - y_i)^2 p(\mathbf{x}_i, y_i) = \frac{1}{2} \sum_{i=1}^N (g(\mathbf{x}_i) - y_i)^2.$$

Now, in the specific case of a linear model like $g(\mathbf{x}, \mathbf{w}) = w_0 + \mathbf{w}^T \mathbf{x}$, it is possible to relate the output through a random error ϵ :

$$y = g(\mathbf{x}, \mathbf{w}) + \epsilon, \quad (12)$$

thus, the error ϵ can be seen graphically as (Figure 4).

Here, a simplifying assumption is that the error comes from a Gaussian Distribution with mean 0 and variance σ^2 , $N(0, \sigma^2)$. Therefore, the Risk Functional can be rewritten as

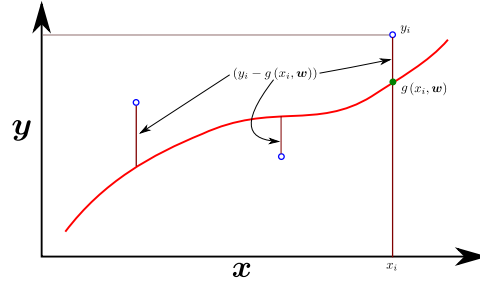


Figure 4: ϵ as the difference $\epsilon = y - g(\mathbf{x}, \mathbf{w})$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \epsilon^2 = \frac{1}{2} \sum_{i=1}^N (w_0 + \mathbf{w}^T \mathbf{x} - y_i)^2 \quad (13)$$

Finally, using our Linear Algebra knowledge, we can rewrite this Risk Functional as

$$E(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (14)$$

with

$$\mathbf{X} = \begin{pmatrix} 1 & (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad (15)$$

The final solution, the estimation $\hat{\mathbf{w}}$, is

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (16)$$

A question arises, how do we use this $\hat{\mathbf{w}}$ to obtain the labels of our classes? We define a threshold saying $(l_1+l_2)/2$ and we define the rule of classification

$$\hat{G}_{LSE}(\mathbf{x}) = \begin{cases} \text{Class 1} & \text{if } \hat{\mathbf{w}}^T \mathbf{x} > (l_1+l_2)/2 \\ \text{Class 2} & \text{if } \hat{\mathbf{w}}^T \mathbf{x} \leq (l_1+l_2)/2 \end{cases}$$

Thus, we have that the hyperplane defined by $\hat{\mathbf{w}}$ is a decision boundary that allows to split the space \mathbb{R}^{d+1} into two classes (Figure 5). As you can see, decision boundary is splitting the space into two sub-planes where, as you can see, you can still have miss-classifications when estimating labels for the data samples.

This represent one of the most popular models in Statistics [11]. Additionally, it is one of the basic models of Learning in Machine Learning.

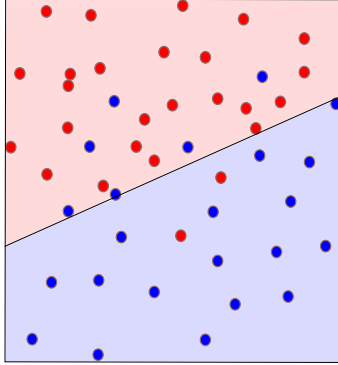


Figure 5: Decision Boundary $\{ \mathbf{x} | \hat{\mathbf{w}}^T \mathbf{x} = (l_1 + l_2)/2 \}$

3.2 Nearest Neighborhood

The other basic model in Machine Learning is the well known k -Nearest Neighbor (k -NN) model. This model can be thought as an answer to the problem of increasing complexity of estimation of h , when the number of dimensions increases beyond a certain level. The rule of classification is quite simple:

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

where $N_K(\mathbf{x})$ is the neighborhood of \mathbf{x} defined by the k closest point \mathbf{x}_i ¹. Here, an important question is: What is the metric used to establish the concept of “near”? For example, we have the concept of near in documents by using the metric:

$$J(D_1, D_2) = \frac{|D_1 \cap D_2|}{|D_1| + |D_2| - |D_1 \cap D_2|} \quad (\text{Jaccard Index})$$

when D_i is the vector that registers if a word appears or not in the document. Another example comes from the idea of using the sample mean and the sample covariance of Data Matrix (Equation 15).

$$\bar{\mathbf{X}} = \frac{1}{N} \sum_{i=1}^N (x_{i1} \quad x_{i2} \quad \cdots \quad x_{ip})^T \quad (\text{Sample Mean}),$$

$$C_{\mathbf{X}} = \frac{1}{N-1} [\mathbf{X} - \bar{\mathbf{X}}]^T [\mathbf{X} - \bar{\mathbf{X}}] \quad (\text{Sample Covariance}).$$

This allows to define the Mahalanobis Distance (Figure 6), and as in the previous Least Squared Error, we can define a new rule of classification based

¹Here, an efficient data structure to find the elements of $N_K(\mathbf{x})$ are known as the $K-d$ trees [12]. This is necessary, if we want to have efficient Learning algorithms for Large Data sets.

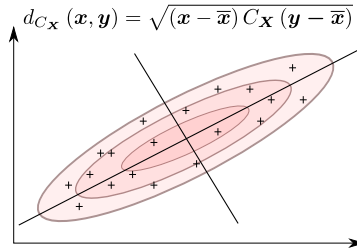


Figure 6: Mahalanobis as a distance measuring the maximum change in each feature.

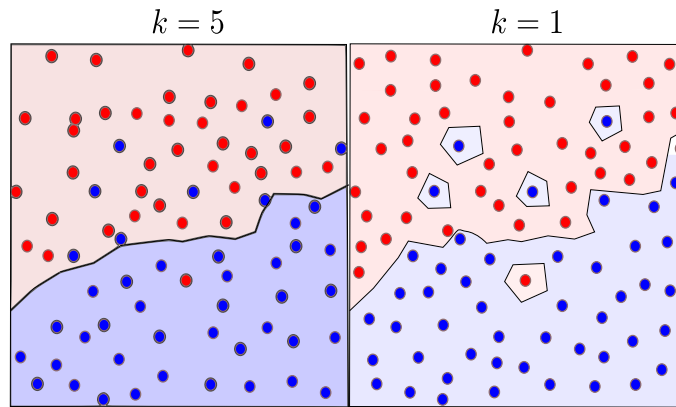


Figure 7: The Decision Boundary $k = 5$ and $k = 1$. Here $k=1$ correspond to a Voronoi diagram [1].

on:

$$\hat{G}_{k\text{-NN}}(\mathbf{x}) = \begin{cases} \text{Class 1} & \text{if } \frac{1}{k} \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} y_i > (l_1 + l_2)/2 \\ \text{Class 2} & \text{if } \frac{1}{k} \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} y_i \leq (l_1 + l_2)/2 \end{cases}$$

where different values for k can produce different boundary decisions (Figure 7).

3.3 Machine Learning Methods as combination and improvements of LSE and k -NN

Something of great interest, as people looking for Patterns in the mathematical structures, mathematical models and data, is the realization that many methods in Machine Learning are related in many ways. In this regards, Hastie et al. [3] have identified two main algorithms for Learning, the Least Squared Error (Section 3.1) and k -Nearest Neighbor (Section 3.2), that can be seen as the two main basis for the algorithms in Machine Learning.

For example, in the case of Support Vector Machines [13], when looking at its dual solution by Vapnik [4]:

$$Q(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \mathbf{x}_j^T \mathbf{x}_i \quad (17)$$

It is known that the elements α_i are different from zero when its associated vector is a support vector. Using an arbitrary kernel in (Equation 17), it is possible to rewrite the equation as:

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{x}_j, \mathbf{x}_i). \quad (18)$$

Thus, if we assume

$$\ell(\alpha_i, \alpha_j, d_i, d_j) = I(\alpha_i > t, \alpha_j > t) \text{sign}(d_i d_j),$$

and $p = 1$, we have

$$\sum_{i=1}^N \left(\alpha_i - \frac{1}{2} \sum_{j=1}^N \ell(\alpha_i, \alpha_j, d_i, d_j) K(\mathbf{x}_j, \mathbf{x}_i) \right)^p \quad (19)$$

Which looks like a Linear Regression with α_i as the desired output and

$$f(\mathbf{x}_i) = \frac{1}{2} \sum_{j=1}^N \ell(\alpha_i, \alpha_j, d_i, d_j) K(\mathbf{x}_j, \mathbf{x}_i)$$

as the estimator of sample \mathbf{x}_i based on their projections by $K(\mathbf{x}_j, *)$.

This is one example of many where algorithms in Machine Learning can be seen as an extension of LSE or k -NN. These could be seen as a little bit limited, but it is useful when trying to analyze and categorize a Machine Learning Algorithm when doing model selection [14].

4 Supervised Learning as Function Approximation

4.1 Regression as Controlled Overfitting

In an exemplifying case, we observe a real-valued input variable $x \in \mathbb{R}$, and we are looking to predict the value of a real valued variable $y \in \mathbb{R}$. Thus, we have a vector of tuples:

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, 2, \dots, N\} \quad (20)$$

for the N samples with their respective outputs. For example, we have the function:

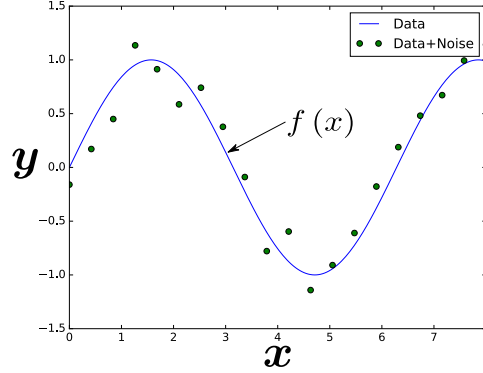


Figure 8: The data and the random observations.

$$g(x) = f(x) + \alpha\epsilon \quad (21)$$

with $\epsilon \sim U(0, 1)$, $\alpha > 0$, $f(x) = \sin\{x\}$ and $x \in \mathbb{R}$ (Figure 8).

Thus, our goal is to exploit this training set to reduce the noise introduced by the system and our observations (Kalman [15]). Basically, we want make predictions of the value $\hat{y} = g(x)$ given a new value \hat{x} . For this is possible to propose the following polynomial:

$$y = g(x, \mathbf{w}) = \sum_{i=0}^d w_i x^i,$$

which can be seen as a generation of new features in a polynomial way (Polynomial Kernels [13]). These functions are linear at the parameter \mathbf{w} and are called linear models. The question, How do we fit the model to our data? Here, we can use the sum of squared errors:

$$R(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N [g(x_i, \mathbf{w}) - y_i]^2 \quad (22)$$

Later in the section of linear models, we will look at ways to get the the canonical solution (Equation 16) for the problem. However, in general, the goal is to obtain a useful approximation (fitting) to $g(\mathbf{x}, \mathbf{w})$ for all \mathbf{x} in some region of \mathbb{R}^d , given the representations in \mathcal{D} . This is an advantage because we can use all the tools generated in the last 200 years for function approximation. Making the task of Approximation/Supervised Learning as **Controlled Overfitting** where in each model a series of parameters are approximated to obtain the desired fitting.

Examples of different approximations that relate the labels and samples in \mathcal{D} are:

1. The linear model $f(x) = \mathbf{x}^T \mathbf{w}$ with $\theta = \mathbf{w}$.
2. The linear basis expansion

$$f_{\theta}(\mathbf{x}) = \sum_{k=1}^K h_k(\mathbf{x}) \theta_k$$

with $h_k(\mathbf{x}) = x_1^2$ or $h_k(\mathbf{x}) = \frac{1}{1+\exp\{-\mathbf{x}^T \theta_k\}}$ and many other ones.

3. The residual sum of squared errors

$$RSS(f_{\mathbf{w}}, \mathbf{x}_0) = \sum_{i=1}^N K_{\lambda}(\mathbf{x}_i, \mathbf{x}_0) (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

Although these are some examples that are quite important in the process of approximation, one should ask the following question, When do we have fittings that are not good at all for the data space? In the following section, we try to answer this question.

4.2 Extreme Fitting, Bias-Variance Problem

For the following development [16, 17], we will take a look at the underlying probability of the samples $\mathbf{x}'s$ by having the following assumptions:

1. There is an underlying distribution for the data $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$ [18]. Here, the notation $p(\mathbf{x}|\Theta)$ represent the random dependencies of the samples of the underlying distribution p which is parametrized by Θ .
2. Some Learning Algorithm has been able to find an estimation function $g(\mathbf{x}|\mathcal{D})$ where the notation $|\mathcal{D}$ represents the dependency of the function g to the training data \mathcal{D} .
3. As in the Least Squared Error, we have seen that the optimal solution for the regression $E_{\mathcal{D}}[y|\mathbf{x}]$ where the subscript \mathcal{D} means that the expected value depends on the distribution $P(X, Y)$.
4. The noise added by the mapping $y = g(x) + \epsilon$ has a distribution $\epsilon \sim N(0, \sigma_{\epsilon}^2)$.

Then, from our probability theory the $Var(X) = E[(X - \mu)^2]$ which can be used to measure the error (Equation 23) of the function $g(\mathbf{x}|\mathcal{D})$ with respect to such expected label of a given sample, $E[y|\mathbf{x}]$.

$$Var_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) = E_{\mathcal{D}} \left(\left(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}||\text{-NN}}[y|\mathbf{x}] \right)^2 \right) \quad (23)$$

Clearly as part of measuring the error, we can use the expected value of the trained function/machine $g(\mathbf{x}|\mathcal{D})$ (Equation 24).

$$E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] \tag{24}$$

the expected output of the machine. Thus, by using a little of algebra, we can rewrite (Equation 23) as

$$\begin{aligned} \text{Var}_{\mathcal{D}} (g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}} \left[(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [y|\mathbf{x}])^2 \right] \\ &= E_{\mathcal{D}} \left[(g(\mathbf{x}|\mathcal{D}) - E[g(\mathbf{x}|\mathcal{D})] + E[g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])^2 \right] \\ &= E_{\mathcal{D}} \left[(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})])^2 + \dots \right. \\ &\quad \dots 2(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}]) + \dots \\ &\quad \left. \dots (E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])^2 \right]. \end{aligned}$$

We can split the terms using the linearity of the expected value:

$$\begin{aligned} \text{Var}_{\mathcal{D}} (g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}} \left[(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})])^2 \right] + \dots \\ &\quad \dots 2E_{\mathcal{D}} [(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])] + \dots \\ &\quad \dots E_{\mathcal{D}} \left[(E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])^2 \right]. \end{aligned}$$

The central term is one of interest

$$E_{\mathcal{D}} [(g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])] = *$$

and by using the linearity of the expected value:

$$\begin{aligned} * &= E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D}) E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}}^2 [g(\mathbf{x}|\mathcal{D})] - \dots \\ &\quad \dots g(\mathbf{x}|\mathcal{D}) E_{\mathcal{D}} [y|\mathbf{x}] + E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] E_{\mathcal{D}} [y|\mathbf{x}]] \\ &= E_{\mathcal{D}}^2 [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}}^2 [g(\mathbf{x}|\mathcal{D})] - \dots \\ &\quad \dots E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] E_{\mathcal{D}} [y|\mathbf{x}] + E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] E_{\mathcal{D}} [y|\mathbf{x}] \\ &= 0. \end{aligned}$$

Finally,

$$\text{Var}_{\mathcal{D}} (g(\mathbf{x}|\mathcal{D})) = \underbrace{E_{\mathcal{D}} \left((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})])^2 \right)}_{\text{VARIANCE}} + \underbrace{(E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])^2}_{\text{BIAS}} \tag{25}$$

Using a little bit of statistics, we can recognize the terms on the $\text{Var}_{\mathcal{D}} (g(\mathbf{x}|\mathcal{D}))$:

1. $E_{\mathcal{D}} \left((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})])^2 \right)$ is **the variance** of output of the learned function $g(\mathbf{x}|\mathcal{D})$.
2. $(E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] - E_{\mathcal{D}} [y|\mathbf{x}])^2$ is a measure on how the expected value of the learned function differs from the expected output, **the bias**.

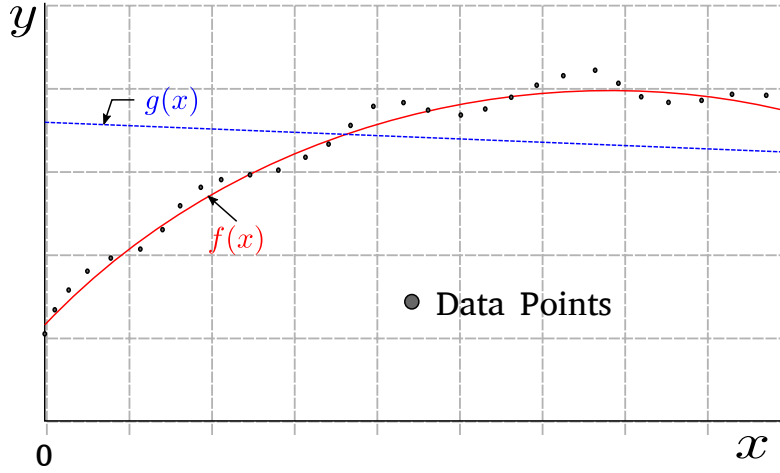


Figure 9: The function used to estimate the output y is highly biased given the mapping.

Thus, we realized that the variance of $g(\mathbf{x}|\mathcal{D})$ ranges between two terms: From variance of the output of the learned function to the bias on the expected output of $g(\mathbf{x}|\mathcal{D})$. It is easy to see this, when X and Y are random variables, then the $g(\mathbf{x}|\mathcal{D})$ is also a random variable. Making of the entire process of learning a random process where we want to avoid extreme cases of fitting.

4.2.1 “Extreme” Cases of Fitting

Case High Bias Imagine that $\mathcal{D} \subset [x_1, x_2]$ in which x lies. For example, we can choose $x_i = x_1 + \frac{x_2 - x_1}{N-1} (i - 1)$ with $i = 1, 2, \dots, N$. Further, we can choose the estimate of $f(x)$, $g(\mathbf{x}|\mathcal{D})$, to be independent of \mathcal{D} . Thus, we can use any function, we can imagine. For example, we could have:

$$g(x) = w_0 + w_1 x$$

Thus, we have the following situation (Figure 9). There, the estimated output is highly biased and far away from the expected outputs of the data samples.

Mathematically, given that $g(x)$ is fixed, we have:

$$E_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] = g(\mathbf{x}|\mathcal{D}) \equiv g(x),$$

with

$$Var_{\mathcal{D}} [g(\mathbf{x}|\mathcal{D})] = 0.$$

On the other hand, because $g(x)$ was chosen arbitrarily the expected bias must be large (Equation 26).

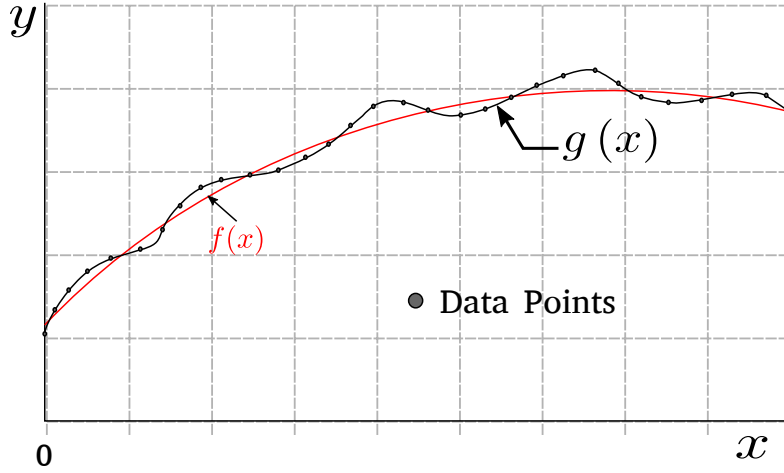


Figure 10: The high variance case of fitting

$$\underbrace{(E_D [g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2}_{BIAS} \quad (26)$$

Clearly, this limited fitting is not what we want, but what could happen if we avoid this high bias and accept a high variance?

Case High Variance In the other hand, $g(x)$ corresponds to a polynomial of high degree so it can pass through each training point in \mathcal{D} (Figure 10).

Now, due to the zero mean of the noise source, we have that:

$$E_D [g(\mathbf{x}|\mathcal{D})] = f(x) = E[y|\mathbf{x}] \text{ for any } x = x_i \quad (27)$$

Thus, at the training points the bias is zero in (Equation 25). However, the variance increases

$$\begin{aligned} E_D \left[(g_1(\mathbf{x}|\mathcal{D}) - E_D [g_1(\mathbf{x}|\mathcal{D})])^2 \right] &= E_D \left[(f(x) + \epsilon - f(x))^2 \right] \\ &= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, \dots, N \end{aligned}$$

In other words, the bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source. Thus, at least in linear models, the Learning procedure needs to make a balance between the Bias and the Variance, this is called the **Bias-Variance trade-off**.

Nevertheless, some observations about this Bias-Variance are due. First, everything that has been said so far applies to both the regression and the classification tasks. However, the Mean Squared Error is not the best way to measure the power of a classifier. This is a classifier that sends everything far away of the hyperplane i.e. away from the values $(+) - 1$.

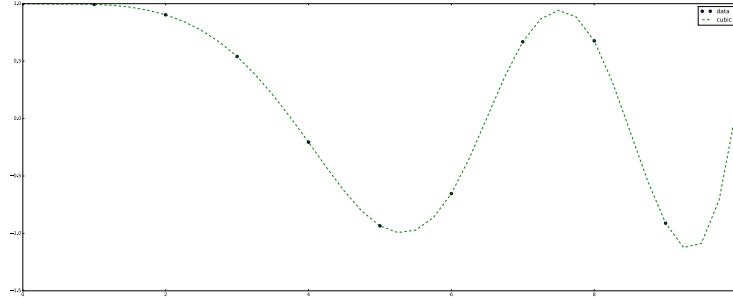


Figure 11: The cubic smooth spline

5 Some Classes of Estimators

Here, Hastie et al. [3] has observed that many of the learning methods fall into different classes of estimators depending on the restrictions imposed by the models. These are only a few ones, but the list is quite larger given the works done in the field, and the need of new algorithms.

5.1 Roughness Penalty

These methods are based on the idea of penalizing the parametrized mapping f to force a regularization over those parameters (Equation 28).

$$R(f; \lambda) = RSS(f) + \lambda J(f) \quad (28)$$

A classic rule to force the regularization is to use a convenient large functional $J(\lambda)$ and a large enough λ to dampen f with high variations over small regions of the input space. For example, the cubic smoothing spline (Figure 11) for one-dimensional input is the solution to the penalized least-squared criterion:

$$R(f; \lambda) = \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda \int \left[\frac{d^2 f(x)}{dx^2} \right]^2 dx \quad (29)$$

Here, the Roughness penalty and the λ controls the high variance of the function f that can be thought as the change in the curvature of the function. Therefore, the penalty works as a way to dampen the high variations in the function f . For example:

1. Linear functions in x , f , appear when $\lambda \rightarrow \infty$.
2. Any function f will stay the same when $\lambda = 0$.

Therefore, penalty or regularization functions express our prior knowledge about the smoothness of the function we are seeking. Thus, it is possible to cast

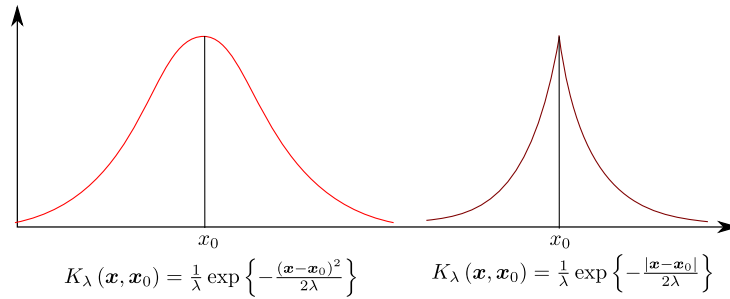


Figure 12: Two types of Kernels

this penalty functions in terms of a Bayesian Framework [19]. The penalty J is the prior distribution, $\sum_{i=1}^N (y_i - f(x_i))^2$ correspond to the likelihood and $PRSS(f; \lambda)$ is the posterior distribution. Thus, minimizing $R(f; \lambda)$ amounts to find the mode in the posterior distribution given that you have a single mode distribution.

5.2 Kernel Methods and Local Regression

You can think on these methods as estimations of regression functions or conditional expectations by specifying:

1. The properties of the local neighborhood.
2. The class of regular functions fitted locally.

For this, they use kernels as

$$K_\lambda(\mathbf{x}, \mathbf{x}_0) = \frac{1}{\lambda} \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}_0\|^2}{2\lambda}\right\}.$$

For example, we have the following kernels (Figure 12). And as in Linear Regression, we can define a way of doing estimation by:

$$R(f_{\mathbf{w}}, \mathbf{x}_0) = \sum_{i=1}^N K_\lambda(\mathbf{x}_i, \mathbf{x}_0) (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2, \quad (30)$$

where $f_{\mathbf{w}}$ could be defined as

1. $f_{\mathbf{w}}(\mathbf{x}) = w_0$ the constant function (Nadaraya-Watson Estimate).
2. $f_{\mathbf{w}}(\mathbf{x}) = \sum_{i=0}^d x_i w_i$ the classic local linear regression models.

In another example the Nearest-Neighbor methods, the following kernel can be defined as:

$$K_k(\mathbf{x}, \mathbf{x}_0) = I[\|\mathbf{x} - \mathbf{x}_0\| \leq \|\mathbf{x}_{(i)} - \mathbf{x}_0\| \mid i = 1, 2, \dots, k]$$

where

1. $\mathbf{x}_{(i)}$ is the training observation ranked i^{th} in distance from \mathbf{x}_0 .
2. $I(S)$ is the indicator of the set S .

5.3 Basis Functions and Dictionary Methods

These are a more wide variety of flexible models defined by the use of linear and polynomial expansions (Equation 31).

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{m=1}^M w_m h_m(\mathbf{x}), \quad (31)$$

where

- h_m is a function on \mathbf{x} .
- with the linear term w_m acting on the function h_m .

Other examples are:

1. Radial basis functions

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{m=1}^M w_m K_{\lambda_m}(\mu_m, \mathbf{x}) \text{ with } K_{\lambda}(\mu, \mathbf{x}) = \exp\left\{-\frac{\|\mathbf{x} - \mu\|^2}{2\lambda}\right\}$$

2. A single-layer feed-forward neural network

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{m=1}^M w_m S(\alpha_m^T \mathbf{x} + b_m) \text{ with } S(y) = \frac{1}{1 + \exp\{-y\}}$$

6 Conclusion

Thus, Machine Learning is a quite wide and vast field with many opportunities that will keep growing as many industries take advantage of this sui generis way of Learning. Nevertheless, this takes time, effort and humility to master because there is no silver bullet that can help to avoid some of the fundamental problems in the field. Only creativity, ingenuity and effort can take the practitioner to calm waters in the vast ocean of Machine Learning.

References

- [1] M. De Berg, M. Van Kreveld, M. Overmars, and O. C. Schwarzkopf, “Computational geometry,” in *Computational geometry*, pp. 1–17, Springer, 2000.
- [2] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [3] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics, Springer New York, 2009.
- [4] S. Haykin, *Neural Networks and Learning Machines*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 2008.
- [5] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*. Academic Press, 1st ed., 2015.
- [6] S. R. Kulkarni, G. Lugosi, and S. S. Venkatesh, “Learning pattern classification—a survey,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2178–2206, 1998.
- [7] R. Silverman, *Modern Calculus and Analytic Geometry*. Dover Books on Mathematics, Dover Publications, 1969.
- [8] C. Robert, *The Bayesian choice: from decision-theoretic foundations to computational implementation*. Springer Science & Business Media, 2007.
- [9] V. Vapnik, *Statistical learning theory*. Wiley, 1998.
- [10] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*, vol. 1. MIT press Cambridge, 2016.
- [11] A. Rencher, *Linear models in statistics*. Wiley series in probability and statistics, Wiley, 2000.
- [12] J. L. Bentley, “Multidimensional binary search trees used for associative searching,” *Commun. ACM*, vol. 18, pp. 509–517, Sept. 1975.
- [13] I. Steinwart and A. Christmann, *Support vector machines*. Springer Science & Business Media, 2008.
- [14] G. Claeskens and N. L. Hjort, *Model Selection and Model Averaging*. No. 9780521852258 in Cambridge Books, Cambridge University Press, April 2008.
- [15] R. E. Kalman, “A new approach to linear filtering and prediction problems,” *Journal of basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.
- [16] P. Domingos, “A unified bias-variance decomposition,” in *Proceedings of 17th International Conference on Machine Learning*, pp. 231–238, 2000.

- [17] S. Geman, E. Bienenstock, and R. Doursat, “Neural networks and the bias/variance dilemma,” *Neural computation*, vol. 4, no. 1, pp. 1–58, 1992.
- [18] R. Ash, *Basic Probability Theory*. Dover Books on Mathematics Series, Dover Publications, Incorporated, 2012.
- [19] M. A. T. Figueiredo, “Adaptive sparseness for supervised learning,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, pp. 1150–1159, Sept. 2003.