
Machine Learning Working Notes

Huang Xiao*

Department of Computer Science
Technical University of Munich
Munich, Germany 85748

Abstract

Machine learning is a fast pacing discipline in many working fields, especially it is now regarded as the most impacting subject in artificial intelligence. In this working notes, I summarize some important notes during my study of machine learning. For the completeness, references are included for readers who are reading this article. Note that this working note is only distributed and shared with author's acknowledge and confirmation. It is not intended as a publishable research paper or tutorial.

1 Linear Models

1.1 Regression

It seems everything starts to grow from linear model, whatever regression or classification, linear models expand to almost many other learning models we face during the research. So starting from a very simple linear regression problem, given training set $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ with n sample, where \mathbf{y} are the responses as numerical values. A linear regression model finds a linear weight vector, which minimizes a certain type of empirical error. This is obviously defined in perspective of statistical learning theory.

$$\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i + b - y_i\|^2 \quad (1)$$

where empirical error introduced by individual sample is equally weighted by $1/n$. It is seen that we are using a straight line to fit a possibly any shaped function, obviously the summation or mean error can be large due to the noise or intrinsic nonlinearity of function. However, a typical misunderstanding of linear model for beginners is that the term linear refers to \mathbf{w} instead of \mathbf{x} . That is to say, we are expecting a linear model on parameters \mathbf{w} , but the feature vectors \mathbf{x} can actually be any shape. Therefore, in literature we mostly see a feature mapping of input \mathbf{x} as $\phi(\mathbf{x})$, and it does not break the linear property of the model. Therefore, we have our linear model as,

$$\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{w}^T \phi(\mathbf{x}_i) + b - y_i\|^2 \quad (2)$$

A typical feature mapping is polynomial feature mapping. Suppose we have a 2-dimensional input data sample (x_1, x_2) , we define a feature mapping as follows,

$$(x_1, x_2) \rightarrow (x_1, x_2, x_1^2 + x_2^2),$$

where a 2-d plane is transformed as a paraboloid in 3-d space. Substituting back into previous linear function, it becomes a polynomial line fitting problem, but it is still linear in \mathbf{w} .

*Website: www.huangxiao.de

To solve the least square problem defined in Eq. 2 to obtain an optimal parameter estimation \mathbf{w} , we take the gradient with respect to the loss and set it to zero. Note that the intercept b can be folded in vector \mathbf{w} by adding additional entry 1 in the end, for simplicity, we ignore it in our formulation. The least square solution is,

$$\mathbf{w}^* = (\Phi(\mathbf{X})^T \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^T \mathbf{Y} \quad (3)$$

As long as the $\Phi(\mathbf{X})^T \Phi(\mathbf{X})$ is not singular, there exists analytical solution. We will call $\Phi(\mathbf{X})$ design matrix, which takes each row as a feature mapping on \mathbf{x} . We will see in short that the inner product of feature mapping can be explicitly defined by a certain kernel function, which established a very important chapter of learning theory, *i.e. Kernel Methods*. To predict the response for a new sample \mathbf{x}^* , we have,

$$\mathbf{y}^* = \phi(\mathbf{x}^*) (\Phi(\mathbf{X})^T \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^T \mathbf{Y}$$

From Probabilistic View

Different from minimizing empirical errors from observations, we can examine the whole problem in a probabilistic view, that is, minimize the uncertainty from observations. If we reformulate the problem as a summation of a deterministic function and a indeterministic noise from a certain probabilistic distribution, we can write the linear model as,

$$\mathbf{y} = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$$

, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ which is defined as a Gaussian noise, the bias term b is again folded in \mathbf{w} . Moreover, we can get the response \mathbf{Y} as a Gaussian distribution as well.

$$\mathbf{Y} \sim \prod_{i=1}^n \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_i), \sigma^2)$$

In order to get the optimal parameter \mathbf{w} , we need to maximize the likelihood $p(\mathbf{Y} | \mathbf{X}, \mathbf{w})$. It equals to maximize the log-likelihood, and the log-likelihood gives,

$$\ln p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) - \frac{1}{2\sigma^2} \|\mathbf{Y} - \Phi(\mathbf{X})\mathbf{w}\|^2 \quad (4)$$

From Eq.4, we can see maximizing the log-likelihood (*MLE*) equals minimizing the least squared error. We will get the exact same solution as in Eq.3.

Overfitting

Now suppose we have four 1-dimensional observations \mathbf{X} , and define an arbitrary feature mapping ϕ , we expect to find a linear model to minimize the least squared error, as we see in (2).

$$\begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix} = \phi \left(\begin{bmatrix} 1.5 \\ 0.5 \\ 2.5 \end{bmatrix} \right) \mathbf{w}$$

If we define the feature mapping ϕ as

$$\begin{bmatrix} 1.5 \\ 0.5 \\ 2.5 \end{bmatrix} \xrightarrow{\phi} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

We can see that \mathbf{w} is exactly the response vector \mathbf{y} , and the least squared error is minimized as zero. A feature mapping can always be defined to achieve zero error, if there's no constraints at all. But obviously, there is no benefit to use this linear model on any prediction task, and mostly likely we will get a very high prediction error based on that. If the model performs well on training dataset, but poorly on unseen data, we can call this situation as overfitting. And certainly, overfitting is a central problem that machine learning attempts to solve.

To avoid overfitting, we can firstly introduce constraint by adding a penalty term on the complexity of the \mathbf{w} , which is known as *regularization*. For example, by penalizing a 2-norm of \mathbf{w} , we can generalize the least squared error problem as,

$$\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{w}^T \phi(\mathbf{x}_i) + b - y_i\|^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (5)$$

Similarly taking the derivative w.r.t. \mathbf{w} and set to zero, we can get the optimal parameters as,

$$\mathbf{w}^* = (\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \lambda \mathbf{I})^{-1} \Phi(\mathbf{X})^T \mathbf{Y} \quad (6)$$

Again from probabilistic view, we introduce uncertainty on \mathbf{w} instead of only considering uncertainty on response \mathbf{y} . Define a prior on \mathbf{w} following a D-dimensional Gaussian distribution,

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$$

We want to capture the posterior distribution on \mathbf{w} after observations of (\mathbf{X}, \mathbf{Y}) , that is, the objective is to maximize the posterior according to Bayes theorem,

$$\max p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \quad (7)$$

The denominator in Eq.7 is also called marginal likelihood, which is independent of \mathbf{w} , therefore, taking the logarithm of the posterior, we have,

$$\ln p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) = -\frac{(n+d)}{2} \ln(2\pi) - n \ln(\Sigma_{\mathbf{w}}) - \frac{1}{2} \ln |\Sigma_{\mathbf{w}}| - \frac{1}{\sigma^2} \|\mathbf{Y} - \Phi(\mathbf{X})\mathbf{w}\|^2 - \frac{1}{2} \mathbf{w}^T \Sigma_{\mathbf{w}}^{-1} \mathbf{w} \quad (8)$$

Take the gradient w.r.t. \mathbf{w} and set to zero, we can get very similar results as in Eq.6.

$$\mathbf{w}_{map} = (\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \sigma^2 \Sigma_{\mathbf{w}}^{-1})^{-1} \Phi(\mathbf{X})^T \mathbf{Y} \quad (9)$$

If the prior is defined with an isotropic Gaussian, we see that the *MAP* solution is equivalent to ℓ_2 regularization form. Now let us look back at the posterior, note that,

$$\begin{aligned} p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) &\propto p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w}) \\ &\propto \exp \left[-\frac{1}{2\sigma^2} (\Phi(\mathbf{X})\mathbf{w} - \mathbf{Y})^T (\Phi(\mathbf{X})\mathbf{w} - \mathbf{Y}) \right] \exp \left(\frac{1}{2} \mathbf{w}^T \Sigma_{\mathbf{w}}^{-1} \mathbf{w} \right) \\ &\propto \exp \left\{ -\frac{1}{2} \mathbf{w}^T \left(\frac{1}{\sigma^2} \Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \Sigma_{\mathbf{w}}^{-1} \right) \mathbf{w} + \frac{1}{\sigma^2} \mathbf{Y}^T \Phi(\mathbf{X}) \mathbf{w} \right\} \end{aligned} \quad (10)$$

By completing the square we can get the mean and covariance of the posterior, that is,

$$\mathbf{w}^* \sim \mathcal{N}(\mathbf{m}^*, \mathbf{A}^{-1})$$

$$\mathbf{m}^* = \frac{1}{\sigma^2} \mathbf{A}^{-1} \Phi(\mathbf{X})^T \mathbf{Y} \quad (11)$$

$$\mathbf{A} = \frac{1}{\sigma^2} \Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \Sigma_{\mathbf{w}}^{-1} \quad (12)$$

And we see that the *MAP* solution is exactly the same as the mode of the posterior. To predict a new input sample \mathbf{x}^* , we can derive a predictive distribution instead of just a single value at the mode, and it is again a Gaussian with posterior mean multiplied with the test sample, and variance of the predictive distribution is the quadratic form on the posterior covariance, which grows with the magnitude of test samples.

$$\begin{aligned} p(\mathbf{y}^* | \mathbf{X}, \mathbf{Y}, \mathbf{x}^*) &= \int p(\mathbf{y}^* | \mathbf{X}, \mathbf{Y}, \mathbf{w}, \mathbf{x}^*) p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) d\mathbf{w} \\ &\sim \mathcal{N}(\Phi(\mathbf{x}^*)^T \mathbf{m}^*, \Phi(\mathbf{x}^*)^T \mathbf{A}^{-1} \Phi(\mathbf{x}^*)) \end{aligned} \quad (13)$$

1.2 Classification

Now let us consider classification problem, where we expect a functional $\pi(\mathbf{x})$ to map input \mathbf{x} to class labels, in binary case $\mathbf{y} = (+1, -1)$. More commonly, we can define a Bernoulli distribution $p(y = +1 | \mathbf{x})$ for one class and $1 - p(y = +1 | \mathbf{x})$ for another. Typically, we would choose a sigmoid function, *e.g.* logistic function or *tanh* to warp a possibly infinite value into a bounding box, *e.g.*, $[0, 1]$ for logistic function. This is a desired behavior, since any function value will be transformed to a probability. A logistic function is defined,

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

Obviously, we also have,

$$\begin{aligned}\sigma(-a) &= 1 - \sigma(a) \\ \tanh(a) &= 2\sigma(2a) - 1 \\ \frac{\partial \sigma(a)}{\partial a} &= (1 - \sigma(a)) \sigma(a)\end{aligned}$$

Thus, the conditional class distribution given input dataset can be defined as a sigmoid function on the linear model $p(y | \mathbf{x}) = \sigma(yf(\mathbf{x}))$, again we fold the bias term b in the parameters \mathbf{w} .

2 Gaussian Process

2.1 Regression

Gaussian process is an important nonparametric regression model which looks for an optimal functional in a space of functions, that minimizes a loss function, although the loss function needs not to be explicitly defined. [4]

Give a training dataset $\mathcal{D} = \{x_i\}_{i=1}^n$ with $x_i \in R^d$, *i.i.d* drawn from certain distribution, we are interested at the predictive distribution of unknown target for the test sample x_* , denoted as f_* . Suppose a prior over \mathbf{y} given input \mathbf{X} is a n -variable Gaussian distribution,

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, K(\mathbf{X}, \mathbf{X}))$$

where $K(\mathbf{X}, \mathbf{X})$ defines a covariance function over \mathbf{X} . Therefore the posterior of f_* given the training dataset \mathcal{D} is also a Gaussian.

$$f_* | \mathbf{y}, \mathbf{X}, x_* \sim \mathcal{N}(\mu_*, \Sigma_*^{-1})$$

where the sufficient statistics can be derived using Bayesian theorem,

$$\bar{f}_* = \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 I)^{-1} \mathbf{y} \quad (14)$$

$$\mathbb{V}(f_*) = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 I)^{-1} \mathbf{k}_* \quad (15)$$

where σ_n^2 is the noise level, \mathbf{K} represents a shorthand for covariance matrix on input \mathbf{X} , and we denote \mathbf{k}_* as the kernel function $k(\mathbf{X}, \mathbf{x}_*)$ for simplicity.

To obtain the Eq.(14)-(15), we can use the following trick. Given two variables (\mathbf{x}, \mathbf{y}) following a Gaussian distribution,

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix} \right)$$

Then we have the conditional distribution,

$$\mathbf{x} | \mathbf{y} \sim \mathcal{N}(\mu_x + CB^{-1}(\mathbf{y} - \mu_y), A - C^T B^{-1} C)$$

Now looking at the predictive mean of training dataset, which can be given by Eq.(14)-(15) on training set \mathbf{X} itself,

$$\bar{\mathbf{f}} = \mathbf{K}(\mathbf{K} + \sigma_n^2 I)^{-1} \mathbf{y}$$

Since \mathbf{K} is symmetric positive definite and its eigendecomposition is $\mathbf{K} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$, where λ_i is the i th eigenvalue and \mathbf{u}_i is the i th eigenvector. Now define a vector $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$, therefore we have $\mathbf{K} = \mathbf{U} \Sigma \mathbf{U}^T$. According the the matrix inverse lemma [[5]], we can derive a simple form for the predictive mean, where we observe that the predictive mean is a linear smooth on their targets.

$$\begin{aligned}
\bar{f} &= \mathbf{K} \left(\sigma_n^{-2} \mathbf{I} - \sigma_n^{-2} \mathbf{I} \mathbf{U} \left(\Sigma^{-1} + \mathbf{U}^T \sigma_n^{-2} \mathbf{I} \mathbf{U} \right)^{-1} \mathbf{U}^T \sigma_n^{-2} \mathbf{I} \right) \mathbf{y} \\
&= \mathbf{K} \left(\sigma_n^{-2} \mathbf{I} - \sigma_n^{-2} \mathbf{U} \left(\Sigma^{-1} + \sigma_n^{-2} \mathbf{I} \right)^{-1} \mathbf{U}^T \sigma_n^{-2} \right) \mathbf{y} \\
&= \left(\sigma_n^{-2} \mathbf{U} \Sigma \mathbf{U}^T - \mathbf{U} \Sigma \begin{bmatrix} \frac{\lambda_1 \sigma_n^{-2}}{\lambda_1 + \sigma_n^2} & & \\ & \ddots & \\ & & \frac{\lambda_n \sigma_n^{-2}}{\lambda_n + \sigma_n^2} \end{bmatrix} \mathbf{U}^T \right) \mathbf{y} \\
&= \left(\sigma_n^{-2} \mathbf{U} \Sigma \mathbf{U}^T - \sigma_n^{-2} \mathbf{U} \begin{bmatrix} \frac{\lambda_1^2}{\lambda_1 + \sigma_n^2} & & \\ & \ddots & \\ & & \frac{\lambda_n^2}{\lambda_n + \sigma_n^2} \end{bmatrix} \mathbf{U}^T \right) \mathbf{y} \\
&= \left(\sigma_n^{-2} \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T - \sigma_n^{-2} \sum_{i=1}^n \frac{\lambda_i^2}{\lambda_i + \sigma_n^2} \mathbf{u}_i \mathbf{u}_i^T \right) \mathbf{y} \\
&= \sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \sigma_n^2} \mathbf{u}_i \mathbf{u}_i^T \mathbf{y} \\
&= \sum_{i=1}^n \frac{\gamma_i \lambda_i}{\lambda_i + \sigma_n^2} \mathbf{u}_i, \quad \text{with } \gamma_i = \mathbf{u}_i^T \mathbf{y}
\end{aligned}$$

2.2 Classification

Using Gaussian process for classification task is a bit more complicated than regression. The main idea of it is to use a ‘squash’ function to convert a predicted function value within $[0, 1]$, *e.g.*, sigmoid function, cumulative Gaussian *pdf*. Given a dataset $\mathcal{D} = \{x_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$ and corresponding labels $\mathbf{y} = [+1, -1]$

3 Support Vector Machines

4 PCA

The basic idea of PCA is to maximally reduce information loss of projecting high dimensional data to lower dimension. Therefore, an intuitive consideration would be that we introduce first of all a projection matrix \mathbf{u} on d -dimensional instance x , so that x is mapped on a lower m -dimensional space. Following column vector routine, we expect that matrix \mathbf{u} as being $m \times d$. Now given a dataset $\mathbf{X} = \{x_i\}_{i=1}^n$, it will be projected on a m -dimensional space by \mathbf{u} . The objective of the projection is to maximize the covariance of data on the lower dimensional space, that is,

$$\max \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}x_i - \mathbf{u}\bar{x}\|^2$$

that can be rewritten as,

$$\begin{aligned}
&\underset{\mathbf{u}}{\text{maximize}} \quad \mathbf{u} \mathbf{S} \mathbf{u}^T \\
&\text{s.t.} \quad \mathbf{u}_i \mathbf{u}_i^T = 1, \quad i = 1, \dots, m
\end{aligned} \tag{16}$$

where \mathbf{S} is the covariance of \mathbf{X} . To prevent \mathbf{u} goes to infinity, we assume \mathbf{u} has unit length, namely, \mathbf{u} represents a set of basis of the lower dimension.

According to (16), we introduce m Lagrangian multipliers $\boldsymbol{\lambda}$ as a diagonal matrix, and we have

$$L = \underset{\mathbf{u}}{\text{maximize}} \quad \mathbf{u} \mathbf{S} \mathbf{u}^T - \boldsymbol{\lambda} \mathbf{u} \mathbf{u}^T$$

Take the derivative of L with respect to \mathbf{u} , we have

$$\begin{aligned} \mathbf{u}S &= \lambda \mathbf{u} \\ S &= \mathbf{u}^{-1} \lambda \mathbf{u} \end{aligned} \quad (17)$$

Since \mathbf{u} is orthogonal, it is therefore not singular. We see that \mathbf{u} and λ are the indeed the eigenvectors and eigenvalues for S respectively. And Eq.(17) is exactly the singular value decomposition of $S = U\Sigma V^T$, we can derive \mathbf{u} and λ from S conveniently.

5 Anomaly Detection: A Survey

In this section, I review literatures of importance in regard of anomaly detection in recent years, a more thorough survey of this research domain can be found also in [3, 1].

6 Supervised Sequence Labeling

The problem of supervised sequence labeling [2] is to assign a sequence of labels given an input sequence. Suppose we have input sequence of length m : $\mathbf{x}_i = \{x_i^0, x_i^1, \dots, x_i^m\}$, we aim to learn a function which assigns a sequence of labels \mathbf{t}_i of length s on \mathbf{x}_i . That is,

$$\begin{aligned} f(\mathbf{x}_i) &= \mathbf{t}_i, \\ \text{where } \mathbf{t}_i &= \{t_i^0, t_i^1, \dots, t_i^s\} \end{aligned}$$

For supervised sequence labeling, there's training labels available and accordingly there're typically three types of learning tasks in sequence labeling.

Sequence Classification It seeks a discriminative classification function to assign singular label to a whole input sequence. This is a m -vs-1 relationship. For instances, a sentence type classifier automatically assigns a sentence to a type label. *Who is president of USA?* is a *question* sentence.

Error function for sequence classification can be as easy as normal binary or multiclass classification problem, e.g., using softmax cross entropy to measure the loss.

Segment Classification Different from sequence classification, segment classification generates a shorter sequence of labels that are assigned to segments of sequence. This is a m -vs- s relationship. Namely we have $|\mathbf{t}_i| < |\mathbf{x}_i|$. This is analog to image segmenation which learns to segment image to different objects.

Error function for segment classification can be set as percentage of misclassified segments. Denote the a test sequence dataset $S' = \{(\mathbf{x}_i, \mathbf{z}_i)_{i=1}^N\}$, and a segment classification function $h(\cdot)$, the segment error can be defined as,

$$E^{seg}(h, S') = \frac{1}{Z} \sum_{(\mathbf{x}, \mathbf{z}) \in S'} \text{HD}(h(\mathbf{x}), \mathbf{z}),$$

where **HD** is the hamming distance between two equal length sequences.

Temporal Classification Temporal classification is a N -vs- N relationship. Each step of input sequence generates a output label, it turns out the output length shall be equal to input sequence length.

A simple error function can be defined as edit distance of two sequences, namely the minimal number of *insertions*, *substitutions* and *deletions*.

$$E^{Temp}(h, S') = \frac{1}{Z} \sum_{(\mathbf{x}, \mathbf{z}) \in S'} \text{ED}(h(\mathbf{x}), \mathbf{z}),$$

where **ED** is edit distance.

References

- [1] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(September):1–58, 2009.
- [2] Alex Graves. *Supervised Sequence Labeling with Recurrent Neural Networks*, volume 12. 2013.
- [3] Victoria J. Hodge and Jim Austin. A Survey of Outlier Detection Methodologies. *Artificial Intelligence Review*, 22(1969):85–126, 2004.
- [4] Carl Edward Rasmussen. *Gaussian processes for machine learning*. MIT Press, 2006.
- [5] Daniel J. Tylavsky and Guy R L Sohie. Generalization of the matrix inversion lemma. *Proceedings of the IEEE*, 74(7):1050–1052, 7 1986.