
Machine Learning Working Notes

Huang Xiao

XIAOHU@IN.TUM.DE

Technische Universität München, Boltzmannstr.3, D-85743 Garching b. München, Germany

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. Gaussian Process

1.1. 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. (Rasmussen, 2006)

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,

$$\mu_* = K_{X_*, X} [K_{X, X} + \sigma_n^2 I]^{-1} \mathbf{y} \quad (1)$$

$$\Sigma_*^{-1} = K_{X_*, X_*} - K_{X_*, X} [K_{X, X} + \sigma_n^2 I]^{-1} K_{X, X_*} \quad (2)$$

Working notes of Machine Learning by Huang Xiao at Technische Universität of München. Copyright 2016 by the author(s).

where σ_n^2 is the noise level and $K_{\cdot, \cdot}$ represents a shorthand for covariance matrix. **Q: Here comes the question of how to estimate the parameters for the covariance K .**

To obtain the Eq.(1)-(2), 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,

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

1.2. 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

2. 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,

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

that can be rewritten as,

$$\begin{aligned} & \underset{\mathbf{u}}{\text{maximize}} \quad \mathbf{u}S\mathbf{u}^T \\ & \text{s.t.} \quad \mathbf{u}_i\mathbf{u}_i^T = 1, \quad i = 1, \dots, m \end{aligned} \quad (3)$$

where 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 (3), we introduce m Lagrangian multipliers λ as a diagonal matrix, and we have

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

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

$$\begin{aligned} \mathbf{u} S &= \lambda \mathbf{u} \\ S &= \mathbf{u}^{-1} \lambda \mathbf{u} \end{aligned} \quad (4)$$

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.(4) is exactly the singular value decomposition of $S = U \Sigma V^T$, we can derive \mathbf{u} and λ from S conveniently.

References

Rasmussen, Carl Edward. Gaussian processes for machine learning. MIT Press, 2006.