

Bayesian regression:

Esmat Hashemi ,Asieh Gholami

s_hashemi1368@yahoo.com,Gholami19@yahoo.com

Department of statistics, payame nor university, Tehran, Iran

Abstract

We are to estimate a point of the W parameter vectors in classic regression models. In contrast, we define uncertainty in W in a Bayesian approach by $p(w)$ probability distribution. The observations from data points change through the Bayesian theorem by the data which are intermediate by probability function in likelihood.

Keyword : Oval-regression- functions

Introduction:

Previous distribution of $p(w)$ indicates uncertainty in w by taking into consideration all data and without losing the generality we can write:

$$1-1: p(w|\alpha) \propto \exp \{ -\alpha \Omega(w) \}$$

We choose special examples. The Gaussian distribution is for $p(w|\alpha)$

As follows:

$$1-2: p(w|\alpha) = \left(\frac{\alpha}{2\pi}\right)^M \exp \left\{ -\frac{\alpha}{2} \|w\|^2 \right\}$$

Now we can use Bayesian theories to state distribution like a priori distribution and posteriori function for w .

Suppose that we want to use previous distribution to find a point estimate for w which is to maximize a posteriori distribution or at least equal to minimize negative logarithm of distribution:

$$1-3: p(w|t, \alpha, \sigma^2) \propto p(w|\alpha)L(w)$$

Where: $L(w) = p(w|t, \alpha, \sigma^2)$

Suppose we want to use a priori distribution to find a point estimate for w which is to maximize a posteriori distribution or equivalently minimize the negative logarithm of distribution, then equal a posteriori distribution logarithm maximization is to minimize.

$$1-4: \quad \frac{1}{2\sigma^2} \sum_{n=1}^N |y(X_n; w) - t_n|^2 + \frac{\alpha}{2} \Omega(w)$$

So we see there is close similarity between such Bayesian and ordinary views on the basis of minimizing error function and the latter was gained by a special approximation on the Bayesian approach; particularly if it is a new value of x , the a priori distribution t is the total sum and having set aside the product from the probability orders on w .

$$1-5: \quad p(t|t, \alpha, \beta) = \int p(w|t, \alpha, \sigma^2) p(t|w, \sigma^2) dw$$

The appropriate variable is unknown in most of the applications though sometimes σ^2 is clear.

1-2: Support vector machine (SVM):

SVM is one of the methods to learn under supervision to be used for stratification and regression; recently it has had a good efficiency compared to previous methods to stratify, for example, perception nervous networks. The SVM function is based on grouping data linearly and in dividing linearly data we try to select a line with more confidence margin. The SVM algorithm solution is grouped as the pattern distinction algorithms. In SVM algorithm whenever it is necessary to distinguish the pattern or group things in special classes it is possible to use it; we provide the pattern matrix and then select the kernel function and then select kernel function parameter and c value in order to use SVM. SVM has been usually defined clearly in such form and leads to prediction based on the function though the application has not been defined as ordinary.

$$1-6: y(x, w) = \sum \omega_i k(x, x_i) + \omega_0$$

Its objective function is to minimize the error measurement in the set and the same time, maximize the margin between two classes. We need $p(t|z)$ distribution for Bayesian activities from logarithm applied programs; such distribution indicates uncertainty in prediction and shows several advantages

such as favorable and flexible decision as an output or other sources of probable data.

1-3: Special models:

We suppose a set of input indicators $\{a_{nd}, t_{an}\}$ where there are a model of turbulence with zero average and σ^2 variance. So:

$$1-7: p(t_n|X) = \mathcal{N}(t_n|y(X_n; w), \sigma^2)$$

Where the $y(x, w)$ function is defined in (1-6). By virtue of tan independence the probability of complete data set may be written as follows:

$$1-8: p(t|w, \sigma^2) = (2\pi\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \|t - \Phi w\|^2\right\}$$

The matrix: $\Phi = [\phi(X_1), \phi(X_2), \dots, \phi(X_N)]^T$ is known as the designer matrix.

$$\phi(X_n) = [1, K(X_n, X_1), K(X_n, X_2), \dots, K(X_n, X_N)]^T$$

And

$$t = (t_1, \dots, t_N)^T \text{ , } w = (w_0, \dots, w_1)^T$$

As many parameters in this model, for example, as the test sample we expect the likelihood maximum method leads to decline to more than the connections. We use Gaussian a priori distribution to complete the characters of previous hierarchy.

$$1-9: p(\omega|\alpha) = \prod_{i=0}^N \mathcal{N}(\omega_i|0, \alpha_i^{-1})$$

These quantities are examples of scale parameter and convenient references were given by gamma for them.

$$1-10: p(\alpha) = \prod_{i=0}^N \text{Gamma}(\alpha_i|a, b)$$

$$1-11: p(\beta) = \text{Gamma}(\beta|c, d)$$

Where: $\text{Gamma}(\alpha|a, b) = \Gamma(a)^{-1} b^a \alpha^{a-1} e^{-b\alpha}$ is gamma function?

Where the gamma reference does not contain useful data in the level of $a \rightarrow 0$ and $b \rightarrow 0$. The ultra parameters may be inconvenient. We may solve this problem by little values: $a=b=c=d=10^{-4}$.

Where:

We consider monotone previous scale is with $a=b=c=d=0$. The Bayesian inference is defined by computations:

$$1-12: p(w, \alpha, \sigma^2 | t) = \frac{p(t|w, \alpha, \sigma^2)p(w, \alpha, \sigma^2)}{p(t)}$$

Then by virtue of a new test point, x_* is the prediction for the goal t_* in previous defined distribution conditions.

$$1-13: p(t_* | t) = \int p(t_* | w, \alpha, \sigma^2 | t) p(w, \alpha, \sigma^2 | t) dw d\alpha d\sigma^2$$

Where we cannot compute directly $p(w, \alpha, \sigma^2 | t)$ from 1-13 because we cannot do ordinary integral.

1-4: Perhaps you ask why Gaussian previous selection should be stated before any priority for dispersed models. We may integrate previous ultra-parameters in order to have such view. It is possible to integrate α independently or take into consideration correctly for previous gammas on additional parameters. It seems previous dispersed like Laplacian reference reached its peak in zero.

$$p(\omega_i) \propto \exp(-|\omega_i|)$$

Which was used to achieve dispersion in Bayesian field.

$$1-14: P(w | t, \alpha) = (2\pi)^{-(N+1)/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(w - \mu)^T \Sigma^{-1}(w - \mu)\right\}$$

This is a posteriori and mean covariance:

1-15:

$$\Sigma = (\sigma^{-2} \Phi^T \Phi + A)^{-1} \quad \mu = \sigma^{-2} \Sigma \Phi^T t$$

$$A = \text{diag}(\alpha_0, \alpha_1, \dots, \alpha)$$

$$\int p(t_* | \alpha, \sigma^2) p(\alpha, \sigma^2 | t) d\alpha d\sigma^2 \simeq p(t_* | \alpha_{MP}, \sigma_{MP}^2)$$

The symbol may be visualization of a mental test when we use the $\Phi_j(x)$ and $\Phi_i(x)$ functions in the same way. Maximum relation:

$$1-16: p(\alpha, \sigma^2 | t) \propto p(t | \alpha, \sigma^2) p(\alpha) p(\sigma^2)$$

Considering α and β for previous ultra-parameter state depends on only logarithm margin probability in $p(t | \alpha, \sigma^2)$ we need following equation to maximize or equalize :

1-17:

$$\begin{aligned} \mathcal{L}(\alpha) &= \ln p(t|\alpha, \sigma^2) = \ln \int_{-\infty}^{\infty} p(t|w, \sigma^2)p(w|\alpha)dw \\ &= -\frac{1}{2} [N \ln 2\pi + \ln|C| + t^T C^{-1}t] \end{aligned}$$

With:

$$1 - 18: C = \sigma^2 I + \Phi A^{-1} \Phi^T$$

1-5: Prediction:

In maximization the prediction is done on the basis of a posteriori distribution by maximizing for a new data x_* .

$$1-19: p(t_*|t, \alpha_{MP}, \sigma_{MP}^2) = \int p(t_*|w, \sigma_{MP}^2)p(w|t, \alpha_{MP}, \sigma_{MP}^2)dw$$

Considering both conditions in the Gaussian may be calculated easily we have:

$$1 - 20: (t_*|t, \alpha_{MP}, \sigma_{MP}^2) = \mathcal{N}(t_*|y_*, \sigma_*^2)$$

$$\text{With: } y_* = \mu^T \phi(x_*)$$

$$\sigma_*^2 = \sigma_{MP}^2 + \phi(x_*)^T \Sigma \phi(x_*)$$

Marginal likelihood:

C analysis is as follows:

1 - 21:

$$C = \sigma^2 I + \sum_{m \neq i} \alpha_m^{-1} \phi_m \phi_m^T + \alpha_i^{-1} \phi_i \phi_i^T = C_{-i} + \alpha_i^{-1} \phi_i \phi_i^T$$

1 - 22:

$$\begin{aligned} \mathcal{L}(\alpha) &= -\frac{1}{2} [N \ln(2\pi) + \ln|C_{-i}| + t^T C_{-i}^{-1}t \\ &\quad - \ln \alpha_i + \ln(\alpha_i + \phi_i^T C_{-i}^{-1} \phi_i) - \frac{(\phi_i^T C_{-i}^{-1}t)^2}{\alpha_i + \phi_i^T C_{-i}^{-1} \phi_i}] \\ &= \mathcal{L}(\alpha_{-1}) + \frac{1}{2} \left[\ln \alpha_i - \ln(\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right] \\ &= \mathcal{L}(\alpha_{-i}) + \ell(\alpha_i) \end{aligned}$$

Where we use following data to simplify:

$$1-23: s_i \triangleq \phi_i^T C_{-i}^{-1} \phi_i \quad , \quad q_i \triangleq \phi_i^T C_{-i}^{-1} t$$

$$\alpha_i = \frac{s_i^2}{q_i^2 - s_i}$$

1-24:

$$\alpha_i = \infty \quad \text{if } q_i^2 \leq s_i$$

Here we begin with the quantity s_i and q_i .

$$1-25: S_i = \phi_i^T C^{-1} \phi_i \quad , \quad Q_i = \phi_i^T C^{-1} t$$

Which is given as follows?

$$(1-26) s_i = \frac{\alpha_i S_i}{\alpha_i - S_i} \quad , \quad q_i = \frac{\alpha_i Q_i}{\alpha_i - S_i}$$

The expressions are when $q_i = Q_i$, $s_i = S_i$ and $\alpha_i = \infty$.

$$1-27: S_i = \phi_i^T B \phi_i - \phi_i^T B \hat{\Phi} \hat{\Sigma} \hat{\Phi}^T B \phi_i$$

$$S_i = \phi_i^T B \phi_i - \phi_i^T B \hat{\Phi} \hat{\Sigma} \hat{\Phi}^T B \phi_i$$

$$\hat{t} \equiv t \text{ and } \hat{\mu} \equiv \mu \text{ and } B \equiv \sigma^{-2} \text{ and } \hat{\Sigma} \equiv \Sigma$$

If Φ_i is in the model, $\alpha_i < \infty$, $q_i^2 \leq S_i$ is still possible after Φ_i .

If Φ_i is not in the model, $\alpha_i = \infty$ and $q_i^2 > S_i$ are possible, Φ_i may be added.

If Φ_i is in the model and $q_i^2 > S_i$, then it is possible to estimate α_i again.

All performances were done to increase the likelihood function; so we have a framework to change the model by adding and eliminating essential function in the probability rule method.

For the variance σ^2 we may we may estimate the equation gained for turbulence variable variance inference.

$$1-28: (\sigma^2)^{new} = \frac{\|t - \Phi\mu\|^2}{N - M + \sum_i \alpha_i \Sigma_{ii}}$$

The dominator N indicates some examples.

1-6: Stratification:

Bayesian stratification is to describe the framework by Bernoulli probability and link function to change in the intended algorithm; as a result, there is an additional approximation step in algorithm.

We use the link function $\sigma(y) = 1/(1 + e^{-y})$ and Bernoulli distribution $P(t|x)$ so we have the density function as follows:

$$1-29: P(t|w) = \prod_{n=1}^N \sigma\{y(X_n; w)\}^{t_n} [1 - \sigma\{y(X_n; w)\}]^{1-t_n}$$

1 – Previous distribution is as follows for present values of α :

$$p(w|t, \alpha) \propto p(t|w) p(w|\alpha)$$

The equation is to find the maximum value:

$$1-30: \ln\{p(t|w)p(w|\alpha)\} =$$

$$\sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] - \frac{1}{2} w^T A w$$

With:

$$y_n = \sigma\{y(X_n; w)\}$$

Second Newton method is to minimize the repeated squares of algorithm to find $\hat{\mu}$.

2- Laplace method is only a second grade approximation.

$$1-31: \nabla_w \nabla_w \ln p(w|t, \alpha) |_{\hat{\mu}} = -(\Phi^T B \Phi + A)$$

$$\text{Where: } B = \text{diag}(\beta_0, \beta_1, \beta_2, \dots, \beta_N)$$

And:

$$\beta_n = \sigma\{y(x_n)\}[1 - \sigma\{y(x_n)\}]$$

In the model $p(w|t, \alpha)$ we may see the local effects of classic problems around μ by:

$$1 - 32: \hat{\Sigma} = (\Phi^T B \Phi + A)^{-1}$$

$$\hat{\mu} = \hat{\Sigma} \Phi^T B \hat{t}$$

1-33:

$$\hat{t} = \Phi\hat{\mu} + B^{-1}(t - \sigma\{\Phi\hat{\mu}\})$$

These equations are used to solve the least squares. We see that the Laplace approximation depends on effectively a regression with indicator t and turbulence data in which reverse variance is presented for n by:

$$\beta_n = \sigma\{y(x_n)\}[1 - \sigma\{y(x_n)\}]$$

1-7: Obvious artificial data: Regression:

The function $\text{sinc}(x) = \sin(x)/X$ is to show the regression vector; we describe the kernel function as follows:

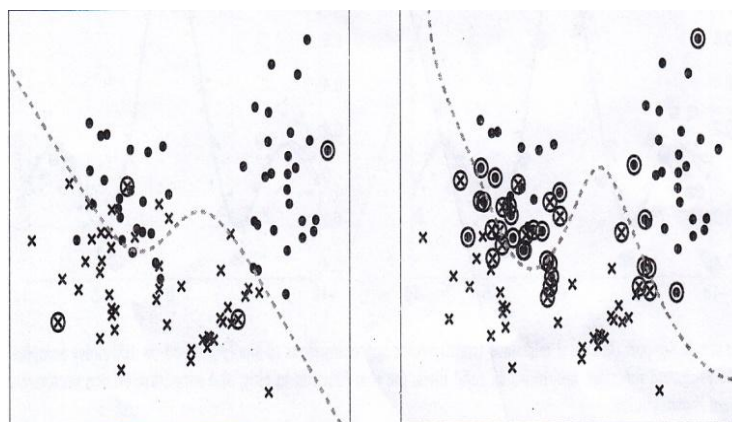
$$1-34: K(x_m, x_n) = 1 + x_m x_n + x_m x_n \min(x_m, x_n) - \frac{x_m + x_n}{2} \min(x_m, x_n)^2 + \frac{\min(x_m, x_n)^3}{3}$$

Where the kernel function to describe a set of primitive function is:

$$\Phi_n(x) = K(x_m, x_n), N=1, \dots, n$$

We use the artificial data produced in 2 later. The class 1 is indicated by x and class two by o in order to show graphically the connection vectors to stratify each two classes. The Bayesian error is about 8 percent by integrating Gaussian 2 mixture in the intervened stratum and the Gaussian kernel is used as follows:

$$1-35: K(X_m, X_n) = \exp(-r^{-2} \|X_m - X_n\|^2)$$



References:

[1] Baser. B. E., Guyton, I. M., and Vapid, V. N. A training algorithm for optimal margin

- [2] Bishop, C.M. and M.E. Tipping (2000). Variational relevance vector machines. In C. Boutilier and M. Goldszmidt (Eds.), Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence, pp. 46-53. Morgan Kaufmann
- [3]oser, B., I. Guyon, and V. N. Vapnik (1992). A training algorithm for optimal margin classifiers. In Proceedings of the Fifth Annual Workshop on Computational Learning Theory, pp. 144-152.
- [4] Classifiers. In D. Haussler, editor, 5th Annual ACM Workshop on COLT, pages 144 –152, Pittsburgh, PA, 1992. ACM Press
- [5] Rich, S. Becker, and Z. Ghahramani (Eds.), Advances in Neural Information Processing Systems 14, pp. 383-389. MIT Press.
- [6] Faull, A. C. and M. E. Tipping (2002). Analysis of sparse Bayesian learning. In T. G.
- [7] Mackay, D. J. C. (1992), The evidence framework applied to classification network. *Neural Computation* 4(5), 720-736.
- [8] Mackay, D. J. C. (1994). Bayesian methods for back propagation network. In E. Doman, J. L. van Hemmen, and K. Schulten (Eds.), *Models of Neural Networks III*, Chapter 6, pp. 211-254. Springer.
- [9] Neal, R. M. (1996). *Bayesian Learning for Neural Networks*. Springer.
- [10] Ripley, B. D. (1996). *Pattern recognition and neural networks*. Cambridge University Press.
- [11] Scholkopf, B., C. J. C. Burges, and A. J. Smola (Eds) (1999). *Advances in Kernel Methods: Support Vector Learning*. MIT Press
- [12] Tipping M. E. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of likelihood maximization for sparse Bayesian models*. In B. Frey and C. M. Bishop (Eds.), *Artificial Intelligence and Statistics*. To appear.
- [13] Tipping, M. E. (2000). The Relevance Vector Machine. In S. A. Sololá, T. K. Len, and K. R. Muller (Eds.), *Advances in Neural Information Processing Systems 12*, pp. 652-658. MIT Press.
- [14] Vapnik, V. N., S. E. Golowich, and A. J. Smola (1997). Support vector method for function approximation, regression estimation and signal processing. In M. Comoser, M. I. Jordan, and T. Pestsche (Eds.), *Advances in Neural Information Processing Systems 9*. MIT Press. *Machine Learning*
- [15] *Bayesian Regression and Classification*.
- [16] Vapid, V. N. (1998). *Statistical Learning Theory*. Wiley.
- [17] Wickham's, P. M. (1995). Bayesian regularization and pruning using a Laplace prior. *Neural Computation* 7(1), 117-143.

[18]Williams, C. K. I. (1997). Prediction with Gaussian processes: From linear regression to liner prediction and beyond. Technical report, NCRG, Aston University, Birmingham, U. K