

Term Project for Machine Learning: Censored Bayesian Polynomial Regression

Mikhail Mishin

April 2017

1 Project Description

Survival analysis is used to analyze data where the outcome variable is the time until occurrence of an event of interest; we call this time as survival time. The event of interest can be, for example, death, outbreak of a disease, divorce, or termination of a computer program.

Data of this kind are often incomplete. Typically, test subjects (e.g., patients) are followed for a fixed time observation period and the time of the occurrence of the event of interest is recorded. However, it may be that in some cases the event of interest did not occur during the period. In these cases, we do not know the exact time of the occurrence but only a lower bound. A datum is called censored if the information of its survival time is incomplete.

In this project we construct a simple Bayesian regression model for censored data. We assume that our data \mathcal{D} consist of N observations. Each observation is a triple (\mathbf{x}_i, y_i, c_i) where $\mathbf{x}_i \in \mathbb{R}^{D'}$ are covariates (predictors, independent variables, input), y_i is the observed survival time and c_i is the censoring time (the end of the observation period). Let z_i be the (unknown) true survival time. Our data are right-censored, that is,

$$y_i = \min(z_i, c_i)$$

We notice that if $y_i < c_i$ we know that $z_i = y_i$. However, if $y_i = c_i$ we know only that $z_i \geq c_i$.

We want to model the true survival times using regression. To this end, we specify a regression model. First, let $\phi(\mathbf{x}_i) : \mathbb{R}^{D'} \rightarrow \mathbb{R}^D$ be a basis function and let $\mathbf{w} \in \mathbb{R}$ be the parameters. Now, we can write

$$z_i = \mathbf{w}^T \phi(\mathbf{x}_i) + \eta_i,$$

where $\eta_i \sim \mathcal{N}(0, \beta^{-1})$ is the residual noise term and β is the noise precision. It follows that $z_i \sim \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$.

Finally, we specify a prior for the parameters \mathbf{w} . We use a Gaussian prior

$$\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I}),$$

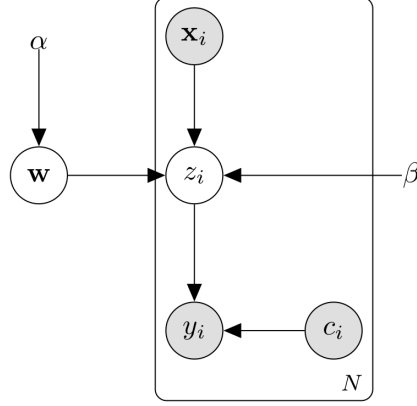


Figure 1: Plate diagram of the censored regression model.

where α is the prior precision.

Ignoring the terms that do not depend on the unknown quantities $\mathbf{z} = \{z_1, z_2, \dots, z_N\}$ and \mathbf{w} , we have the following joint distribution

$$p(\mathbf{w}, \mathbf{z}, \mathcal{D}, \alpha, \beta) \propto p(\mathbf{w}|\alpha) \prod_{i=1}^N p(z_i|\mathbf{w}, \mathbf{x}_i, \beta) p(y_i|z_i, c_i).$$

This distribution is illustrated in Figure 1. We use hyperparameter values $\alpha = \beta = 1$ and omit them from conditional distributions for simplicity.

In this project, we consider polynomial regression. We assume that $\mathbf{x}_i = x_i$ is a singleton, that is, $D' = 1$. To obtain d :th degree polynomial regression, we specify a basis function $\phi^{(d)}(x) = [1, x^1, x^2, \dots, x^d]^T$. The constant term corresponds to the intercept. As a result, we have $D = d + 1$.

Due to the latent variables z_i , it is not possible to obtain the posterior distribution in analytic form here. Fortunately, it is possible to use the EM algorithm to obtain a maximum a posteriori (MAP) estimate for the parameters \mathbf{w} . In Bayesian EM, we want to find

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \log p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} \log p(\mathbf{w}, \mathcal{D})$$

The task for this project is to derive, implement, and test an EM algorithm for censored Bayesian polynomial regression.

2 Derivations

First, we derive the log-posterior functions for complete data and incomplete data. Then, we derive the EM update equations for the model.

2.1 Complete data log-posterior

We can treat data \mathcal{D} as constant. Thus, the posterior is proportional to the joint distribution. The joint distribution function for the complete data is given by:

$$p(\mathbf{w}, \mathbf{z}|\mathcal{D}) \propto p(\mathbf{w}, \mathbf{z}, \mathcal{D}) = p(\mathbf{w}) \prod_{i=1}^N p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i)$$

Taking the logarithm we obtain:

$$\log p(\mathbf{w}, \mathbf{z}|\mathcal{D}) = \log p(\mathbf{w}) + \sum_{i=1}^N [\log p(z_i|\mathbf{w}, \mathbf{x}_i) + \log p(y_i|z_i, c_i)] + \text{const} \quad (1)$$

2.2 Incomplete data log-posterior

The log-posterior function for the incomplete data is given by:

$$\begin{aligned} \log p(\mathbf{w}|\mathcal{D}) &= \log \int_{-\infty}^{\infty} p(\mathbf{w}, \mathbf{z}|\mathcal{D}) d\mathbf{z} \\ &\propto \log \int_{-\infty}^{\infty} p(\mathbf{w}, \mathbf{z}, \mathcal{D}) d\mathbf{z} \\ &= \log \int_{-\infty}^{\infty} p(\mathbf{w}) \prod_{i=1}^N p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i) d\mathbf{z} \\ &= \log p(\mathbf{w}) + \log \int_{-\infty}^{\infty} \prod_{i=1}^N p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i) d\mathbf{z} \end{aligned}$$

Moving the product out of the integral

$$\begin{aligned} &= \log p(\mathbf{w}) + \log \prod_{i=1}^N \int_{-\infty}^{\infty} p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i) dz_i \\ &= \log p(\mathbf{w}) + \sum_{i=1}^N \log \int_{-\infty}^{\infty} p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i) dz_i \end{aligned} \quad (2)$$

where

$$\log p(\mathbf{w}) = -\frac{1}{2} \log |2\pi\alpha^{-1}\mathbf{I}| - \frac{1}{2} \mathbf{w}^T \alpha \mathbf{I} \mathbf{w} \quad (3)$$

Observe that $p(y_i|z_i, c_i)$ is a delta distribution. Consider two cases:

- $y_i < c_i$:

$$p(y_i|z_i, c_i) = \begin{cases} 1, & \text{if } z_i = y_i, \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

- $y_i = c_i$:

$$p(y_i|z_i, c_i) = \begin{cases} 1, & \text{if } z_i \geq c_i, \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

Therefore, the integral in (2) evaluates to:

$$\int_{-\infty}^{\infty} p(z_i|\mathbf{w}, \mathbf{x}_i) p(y_i|z_i, c_i) dz_i = \begin{cases} p(Z_i = y_i|\mathbf{w}, \mathbf{x}_i), & \text{if } y_i < c_i, \\ \int_{c_i}^{\infty} p(z_i|\mathbf{w}, \mathbf{x}_i) dz_i, & \text{if } y_i = c_i \end{cases} \quad (6)$$

We can write the final formula for the incomplete data log-posterior using an indicator function:

$$\begin{aligned} \log p(\mathbf{w}|\mathcal{D}) = \log p(\mathbf{w}) + \sum_{i=1}^N [& (1 - \mathcal{I}(y_i = c_i)) \log f(y_i|\mathbf{w}^T \mathbf{x}_i, \beta^{-1}) \\ & + \mathcal{I}(y_i = c_i) \log(1 - F(c_i|\mathbf{w}^T \mathbf{x}_i, \beta^{-1}))] \end{aligned} \quad (7)$$

where $f(u|\mu, \sigma^2)$ and $F(u|\mu, \sigma^2)$ are the PDF and CDF of the Gaussian distribution with the mean μ and variance σ^2 , respectively.

2.3 EM Update Equations

In Bayesian EM, we want to find $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \log p(\mathbf{w}, \mathcal{D})$.

- **E-step** To update z_i consider two cases:

- If $y_i < c_i$,

$$\left\langle z_i \right\rangle_{Z|\mathcal{D}, \mathbf{w}_t} = \left\langle z_i | z_i = y_i, y_i \right\rangle = y_i \quad (8)$$

- If $y_i = c_i$,

$$\left\langle z_i \right\rangle_{Z|\mathcal{D}, \mathbf{w}_t} = \left\langle z_i | z_i \geq c_i, \mathbf{w}_t, \mathbf{x}_i \right\rangle = \mathbf{w}_t^T \phi(\mathbf{x}_i) + \beta^{-\frac{1}{2}} \text{H} \left(\frac{c_i - \mathbf{w}_t^T \phi(\mathbf{x}_i)}{\beta^{-\frac{1}{2}}} \right) \quad (9)$$

, where

$$\text{H}(u) = \frac{f(u)}{1 - F(u)}$$

The last equation follows from the formula for the expectation of a normally distributed random variable bounded below ($z_i \geq c_i$).

- **M-step** Expected complete data joint distribution:

$$\begin{aligned}
R(\mathbf{w}, \mathbf{w}_t) &= \left\langle \log p(\mathbf{w}, Z, \mathcal{D}) \right\rangle_{Z|\mathcal{D}, \mathbf{w}_t} \\
&= \log p(\mathbf{w}) + \sum_{i=1}^N \left\langle \log p(z_i | \mathbf{w}, \mathbf{x}_i) \right\rangle_{Z|\mathcal{D}, \mathbf{w}_t} + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \alpha \mathbf{I} \mathbf{w} - \sum_{i=1}^N \left\langle \frac{(z_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2}{2\beta^{-1}} \right\rangle_{Z|\mathcal{D}, \mathbf{w}_t} + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \alpha \mathbf{I} \mathbf{w} - \sum_{i=1}^N \frac{\langle z_i^2 \rangle - 2\langle z_i \rangle \mathbf{w}^T \phi(\mathbf{x}_i) + (\mathbf{w}^T \phi(\mathbf{x}_i))^2}{2\beta^{-1}} + \text{const}
\end{aligned} \tag{10}$$

Let us take the derivative of $R(\mathbf{w}, \mathbf{w}_t)$ with respect to \mathbf{w} :

$$\begin{aligned}
\frac{\partial R(\mathbf{w}, \mathbf{w}_t)}{\partial \mathbf{w}} &= -\alpha \mathbf{I} \mathbf{w} + \frac{\beta}{2} \sum_{i=1}^N 2\langle z_i \rangle \phi(\mathbf{x}_i) - 2(\mathbf{w}^T \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i) \\
&= -\alpha \mathbf{I} \mathbf{w} + \beta \sum_{i=1}^N \langle z_i \rangle \phi(\mathbf{x}_i) - \beta \left(\sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \right) \mathbf{w} \\
&= - \left(\beta \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T + \alpha \mathbf{I} \right) \mathbf{w} + \beta \sum_{i=1}^N \langle z_i \rangle \phi(\mathbf{x}_i) \tag{11}
\end{aligned}$$

Maximizing for \mathbf{w} :

$$\frac{\partial R(\mathbf{w}, \mathbf{w}_t)}{\partial \mathbf{w}} = 0 \implies \mathbf{w} = \left(\beta \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T + \alpha \mathbf{I} \right)^{-1} \beta \sum_{i=1}^N \langle z_i \rangle \phi(\mathbf{x}_i) \tag{12}$$

3 Evaluation

We estimate MAP parameters for polynomials with degrees $d = 1, \dots, 5$ in Table 1. In our implementation of the EM algorithm we use 30 iterations and 10 random restarts. We plot the MAP curves using the estimated parameters in Figure 2. For comparison, we also plot the MAP curves for the standard regression for the model $y_i = \mathbf{w}^T \phi(\mathbf{x}_i) + \eta_i$. We use the analytic formula to find the parameters (omitting the derivation):

$$\mathbf{w} = \left(\beta \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T + \alpha \mathbf{I} \right)^{-1} \beta \sum_{i=1}^N y_i \phi(\mathbf{x}_i) \quad (13)$$

Comparing two curves, one can see how the censored regression takes into account the censored data that is ignored by the standard regression. To verify the convergence of our algorithm we also plot incomplete data log-posterior as a function of iteration (Figure 3). Indeed, log-posterior increases after each M-step and converges quickly.

Degree	Estimated MAP parameters
1	[5.2797; 1.2552]
2	[5.7734; 1.0746; -0.4908]
3	[5.8030; 1.2196; -0.5033; -0.0582]
4	[5.7066; 1.2963; -0.1504; -0.0794; -0.1048]
5	[5.7278; 1.4491; -0.1922; -0.2818; -0.0958; 0.0468]

Table 1: Estimated MAP parameters for censored regression

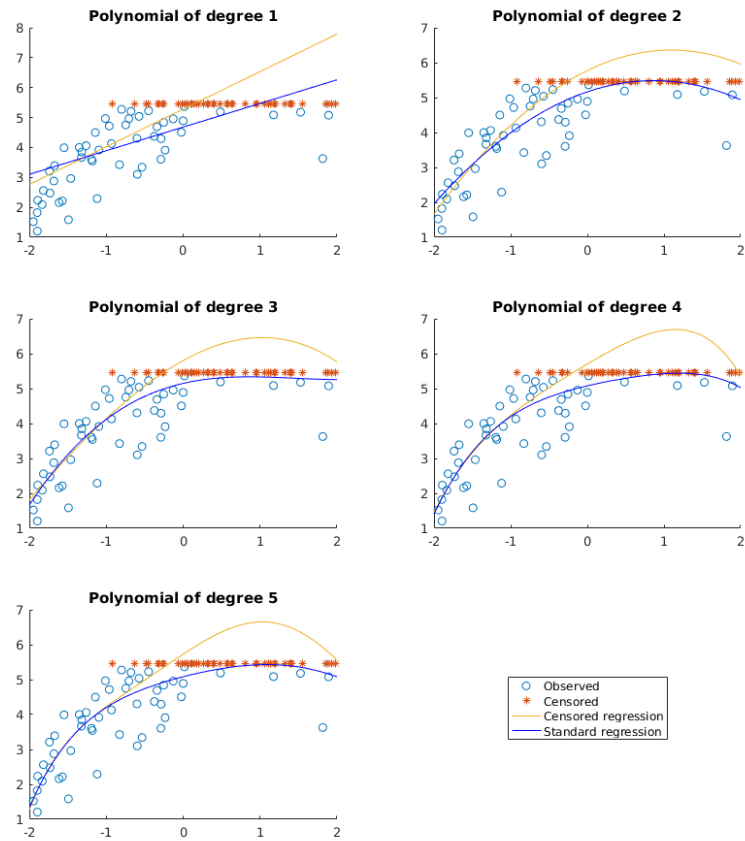


Figure 2: MAP curves for censored and standard regression

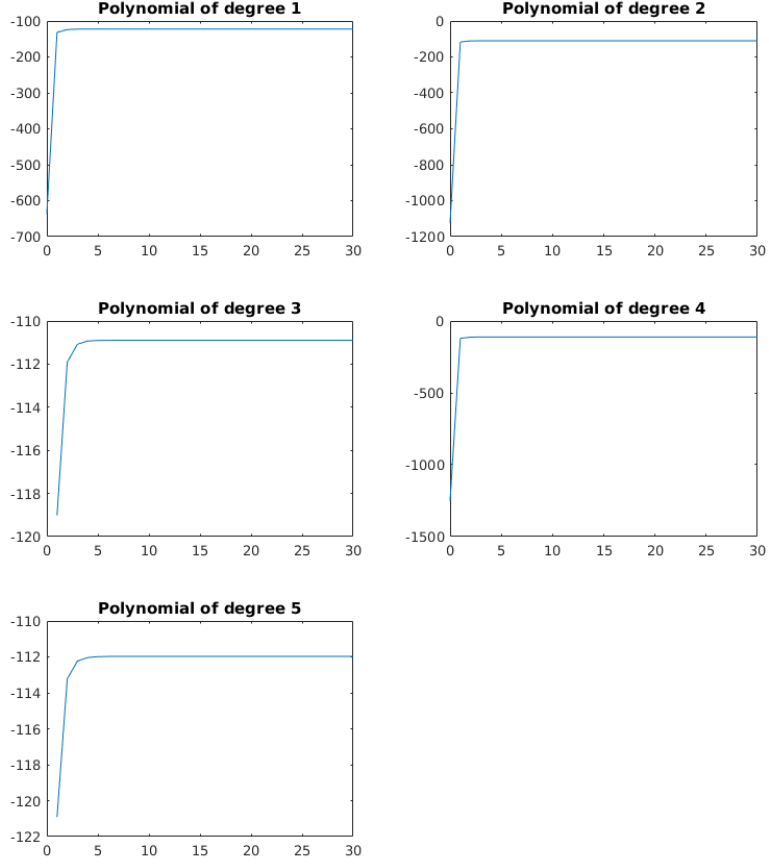


Figure 3: Incomplete data log-posterior as a function of iteration

4 Model Selection

To select a model we compute Bayesian information criterion (BIC) for each d (Figure 4). We choose the model with the lowest BIC - polynomial of degree $d = 2$. Looking at the MAP curves from Figure 2 again, we can say that our choice is consistent and reasonable.

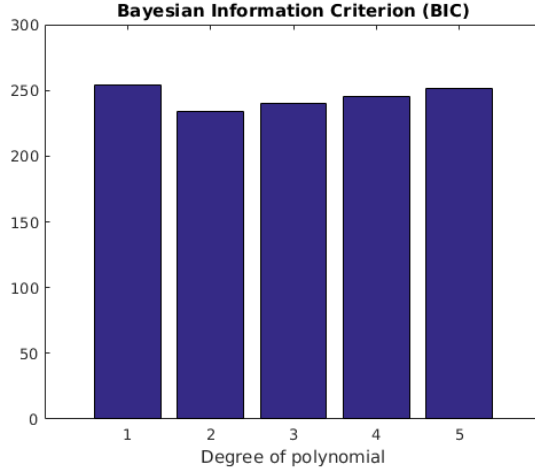


Figure 4: Bayesian information criterion (BIC)

5 Implementation

Code for computing the log-posterior function for the incomplete data:

```

1 function [ logpst ] = logPosterior( phi_x, y, c, w, alpha, beta )
2 %LOGPOSTERIOR Compute incomplete data log-posterior
3     d = length(w);
4     log_pw = -d/2*log(2*pi*1/alpha) - 1/2*(w'*alpha*eye(d)*w);
5     sum_log_pz = 0;
6     mean_z = phi_x*w;
7     for i = 1:length(y)
8         if c(i)
9             % y_i = c_i
10            sum_log_pz = sum_log_pz + ...
11                log(1 - normcdf(y(i), mean_z(i), beta^(-1/2)));
12        else
13            % y_i < c_i
14            sum_log_pz = sum_log_pz + ...
15                log(normpdf(y(i), mean_z(i), beta^(-1/2)));
16        end
17    end
18    logpst = log_pw + sum_log_pz;
19 end

```

Implementation of the EM algorithm:

```

1 for d = 1:max_degree
2
3     best_log_post = -Inf;
4     best_w = -Inf*ones(d+1, 1); % MAP estimate for parameters
5     phi_x = polyBasis(x, d);
6     S = inv(beta*(phi_x'*phi_x) + alpha*eye(d+1));
7     E_z = zeros(size(y));
8

```

```

9   for s = 1:n_starts
10      % Initialize weights
11      w = mvnrnd(zeros(1,d+1), alpha^(-1/2)*eye(d+1))';
12      % incomplete data log-posterior
13      log_posts = zeros(n_iter + 1, 1);
14      % log-posterior before first iteration
15      log_posts(1) = logPosterior(phi_x, y, c, w, alpha, beta);
16
17      for i = 1:n_iter
18         % E-step
19         mean_z = phi_x*w;
20         for j = 1:n_obs
21            if c(j)
22               alpha2 = (y(j) - mean_z(j))/beta^(-1/2);
23               E_z(j) = mean_z(j) + beta^(-1/2)*H_function(
alpha2);
24            else
25               E_z(j) = y(j);
26            end
27         end
28         % M-step
29         w = beta*S*phi_x'E_z;
30
31         % Compute log-posterior
32         log_posts(i+1) = logPosterior(phi_x, y, c, w, alpha,
beta);
33      end
34      if log_posts(n_iter+1) > best_log_post
35         best_log_post = log_posts(n_iter+1);
36         best_w = w;
37         log_postconvergences = log_posts;
38      end
39
40   end
41
42   % Compute BIC
43   bics(d) = bic(best_log_post, d, n_obs);
44
45   % Store MAP estimate
46   ws{d} = best_w;
47 end

```

Code for computing polynomial basis functions:

```

1 function phi_x = polyBasis(x, d)
2 %POLYBASIS Polynomial basis function
3   phi_x = bsxfun(@power, x, 0:d);
4 end

```

Code for computing Bayesian information criterion:

```

1 function [ crt ] = bic( log_post, d, n_obs )
2 %BIC Bayesian information criterion
3   crt = (d+1)*log(n_obs) - 2*log_post;
4 end

```

6 Acknowledgements

The section Project Description is taken from the original description for the course written by Pekka Marttinen, Pekka Parviainen, Pedram Daei, and Sami Remes.