

Assignment Project Exam Help

Predictive Analytics

Week 6: Model Selection and Estimation III

<https://powcoder.com>

Semester 2, 2018

Discipline of Business Analytics, The University of Sydney Business School

Add WeChat powcoder

Week 6: Model Selection and Estimation III

1. Maximum likelihood (continued)

2. Inference for the ML estimator (optional)

3. Analytical criteria

4. Comparison of model selection methods

5. Limitations of model selection

6. Optimism (optional)

Reading: Chapter 6.1 of ISL.

Exercise questions: Chapter 6.1 of ISL, Q1. Try to answer this question based on your existing knowledge of regression variable selection, which will help you be prepared for the next lecture.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Maximum likelihood (continued)
<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Maximum likelihood estimation (MLE), which we have discussed in the context of linear regression is one of the most important concepts in statistics. We now present it more generally for inference.

<https://powcoder.com>
Add WeChat powcoder

ML for discrete distributions (key concept)

Let $p(y; \theta)$ be a discrete probability distribution. The likelihood function is

Assignment Project Exam Help

$$\begin{aligned}\ell(\theta) &= P(Y_1 = y_1, Y_2 = y_2, \dots, Y_N = y_N) \\ &= P(Y_1 = y_1) P(Y_2 = y_2) \dots P(Y_N = y_N)\end{aligned}$$

$$= \prod_{i=1}^N p(y_i; \theta)$$

Add WeChat powcoder

The maximum likelihood estimate $\hat{\theta}$ is the value of θ that maximises $\ell(\theta)$.

ML for continuous distributions (key concept)

Let $p(y; \theta)$ be a density function. The likelihood function is

Assignment Project Exam Help

$$\begin{aligned}\ell(\theta) &= p(y_1; \theta) p(y_2; \theta) \cdots p(y_N; \theta) \\ &= \prod_{i=1}^N p(y_i; \theta)\end{aligned}$$

<https://powcoder.com>

Add WeChat powcoder

The maximum likelihood estimate $\hat{\theta}$ is the value of θ that maximises $\ell(\theta)$.

Assignment Project Exam Help

- Even though $\ell(\theta)$ equals an expression that involves $p(y_i; \theta)$, we think of these functions in different ways.

<https://powcoder.com>

- When considering a probability mass function or density $p(y; \theta)$, we consider y to be a variable, and θ to be fixed.

Add WeChat powcoder

- In the likelihood, θ is a variable, and y is fixed.

Log-likelihood (key concept)

The log-likelihood is

Assignment Project Exam Help

<https://powcoder.com>

$$L(\theta) = \log \left(\prod_{i=1}^N p(y_i; \theta) \right)$$
$$= \sum_{i=1}^N \log p(y_i; \theta)$$

Add WeChat powcoder

Because the log-likelihood is a **monotonic** transformation of the likelihood, maximising it is the same as maximising the likelihood.

Example: Bernoulli distribution

Suppose that Y_1, \dots, Y_N follow the Bernoulli distribution with parameter θ (the probability of a success).

Assignment Project Exam Help

$$p(y_i; \theta) = \theta^{y_i} (1 - \theta)^{(1-y_i)}$$

<https://powcoder.com>

$$\ell(\theta) = \prod_{i=1}^N \theta^{y_i} (1 - \theta)^{(1-y_i)}$$

Add WeChat powcoder

$$\begin{aligned} L(\theta) &= \sum_{i=1}^N [y_i \log(\theta) + (1 - y_i) \log(1 - \theta)] \\ &= \left(\sum y_i \right) \log(\theta) + (N - \sum y_i) \log(1 - \theta) \end{aligned}$$

Example: Bernoulli distribution

Derivative of the log-likelihood with respect to θ :

$$\frac{dL(\theta)}{d\theta} = \frac{\sum y_i}{\theta} - \frac{N - \sum y_i}{1 - \theta}$$

The ML estimate therefore satisfies

$$\frac{\sum y_i}{\hat{\theta}} = \frac{N - \sum y_i}{1 - \hat{\theta}}.$$

The solution is the sample proportion:

$$\hat{\theta} = \frac{\sum_{i=1}^n y_i}{N}.$$

What about the 2nd order derivative?

Assignment Project Exam Help

Inference for the ML estimator

(optional) <https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

The score function is

$$s(\hat{\theta}) = \nabla L(\theta)|_{\theta=\hat{\theta}}$$

<https://powcoder.com>

For example, when the parameter is a scalar

$$s(\hat{\theta}) = \sum_{i=1}^n \frac{d \log p(y_i; \theta)}{d\theta} \Big|_{\theta=\hat{\theta}}$$

Inference for the ML estimator

The **observed information matrix** is the negative of the second derivative (the Hessian matrix) of the log-likelihood.

$$J(\hat{\theta}(\mathcal{D})) = -\nabla_{\theta}^2 L(\theta) \Big|_{\theta=\hat{\theta}}$$

<https://powcoder.com>

When the parameter is a scalar,

$$J(\hat{\theta}) = \sum_{i=1}^N \frac{d^2 \log p(y_i)}{d\theta^2} \Big|_{\theta=\hat{\theta}}.$$

Add WeChat powcoder

Assignment Project Exam Help

We define the **Fisher information matrix** as the expected value of the observed information matrix

<https://powcoder.com>

$$I_n(\hat{\theta}) = E_{\theta} [J(\hat{\theta}(\mathcal{D}))]$$

So the **observed information matrix** is a sample based version of the Fisher information matrix.

Add WeChat powcoder

Inference for the ML estimator

Assignment Project Exam Help
A standard result shows that the sampling distribution of the ML estimator converges to the normal distribution

$$\hat{\theta} \rightarrow N(\theta, I_n^{-1}(\theta))$$

as $n \rightarrow \infty$.

That suggests the large sample approximations

$$N(\theta, I_N(\hat{\theta})^{-1}) \text{ or } N(\theta, J(\hat{\theta})^{-1})$$

Example: Bernoulli distribution

Continuing the example, the observed information matrix is

$$\frac{d^2 L(\theta)}{d\theta^2} = -\frac{\sum y_i}{\theta^2} - \frac{n - \sum y_i}{(1-\theta)^2}$$

Assignment Project Exam Help

<https://powcoder.com>

Since $E(Y) = \theta$,

$$E(J(\theta)) = \frac{N}{\theta(1-\theta)},$$

so that

$$I_N^{-1} = \frac{\theta(1-\theta)}{N},$$

which is familiar as the variance of a sample proportion from basic statistics.

Add WeChat powcoder

Inference for the ML estimator

The corresponding estimates for the standard errors of individual parameters are

$$SE(\hat{\theta}_j) = \sqrt{I_{jj}(\hat{\theta})^{-1}} \text{ or } SE(\hat{\theta}_j) = \sqrt{I(\hat{\theta})_{jj}^{-1}}$$

A large sample $100 \times (1 - \alpha)\%$ confidence interval is

$$\hat{\theta}_j \pm z_{\alpha/2} \times SE(\hat{\theta}_j)$$

Assignment Project Exam Help

The following large sample approximation leads to accurate confidence intervals and hypothesis tests

<https://powcoder.com>

$$2 \left(L(\hat{\theta}) - L(\theta) \right) \sim \chi_d^2,$$

where d is the number of parameters in θ .

Add WeChat powcoder

Assignment Project Exam Help

Analytical criteria
<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Analytic criteria provide estimate the generalisation error based on theoretical arguments. They have the form:

<https://powcoder.com>

$\text{criterion} = \text{training loss} + \text{penalty for number of parameters}$

Add WeChat powcoder

Mallow's C_p statistic

The **Mallow's** C_p statistic applies to linear regression. It directly implements the recipe suggested by our calculation of the optimism:

<https://powcoder.com>

$$C_p = \frac{\text{RSS}}{N} + \frac{2}{N} \hat{\sigma}^2 (p + 1),$$

Add WeChat powcoder

In the formula, $\hat{\sigma}^2$ is an estimate of variance of the errors based on the largest model under consideration.

Assignment Project Exam Help

We select the model with the lowest C_p . To compare two specifications,

<https://powcoder.com>

$$\Delta C_p = \text{MSE}_1 - \text{MSE}_2 + \frac{2}{N} \hat{\sigma}^2 (p_1 - p_2).$$

Add WeChat powcoder

Akaike Information Criterion (key concept)

Assignment Project Exam Help

The **Akaike information criterion (AIC)** applies to models estimated by maximum likelihood.

<https://powcoder.com>

$$\text{AIC} = -2L(\hat{\theta}) + 2p,$$

where $L(\hat{\theta})$ is the maximised log-likelihood and p is the number of estimated parameters. We select the model with the lowest AIC.

Add WeChat powcoder

Assignment Project Exam Help

- The AIC is one of the most popular and versatile strategies for model selection.
- The formula follows the in-sample performance plus penalty for complexity structure.
- The AIC has a rigorous theoretical justification which we not address here. However, keep in mind that it is an asymptotic approximation ($N \rightarrow \infty$).

<https://powcoder.com>

Add WeChat powcoder

AIC for linear regression

In the special case of comparing linear regression specifications

under Gaussian errors, the AIC simplifies to (up to proportionality and ignoring constant terms in the log-likelihood):

$$\text{AIC} \propto \log \left(\frac{\text{RSS}}{N} \right) + \frac{2}{N} (p + 2).$$

The number of parameters is $\ell = p + 2$ because the parameter vector includes the constant and the variance of errors. Note that this is different from the formula in the book, which is a simplification with unknown error.

Relation between Mallows's C_p and the AIC

For a linear regression with Gaussian errors and known variance:

Assignment Project Exam Help

$$AIC = \frac{1}{\hat{\sigma}^2} \left(\frac{RSS}{N} + \frac{2}{N} \hat{\sigma}^2 (p+1) \right),$$

which compares to

<https://powcoder.com>

$$C_p = \frac{RSS}{N} + \frac{2}{N} \hat{\sigma}^2 (p+1).$$

Add WeChat powcoder

Hence, the AIC and C_p lead to the same decision in this case. For practical purposes, the AIC and C_p are regarded as the same for linear regression.

Bayesian information criterion

Assignment Project Exam Help

The Bayesian information criterion (BIC) also applies to models estimated by maximum likelihood.

$$\text{BIC} = -2L(\hat{\theta}) + \log(N)p$$

where $L(\hat{\theta})$ is the maximised log-likelihood, p is the number of estimated parameters, and N is the sample size. We select the model with the lowest BIC.

Add WeChat powcoder

Assignment Project Exam Help

- The BIC formula is comparable to the AIC case, but with 2 penalty factor replaced by $\log(N)$. Hence, the BIC penalises complexity more heavily when $N \geq 8$. The BIC has a very different theoretical justification to the AIC.

- The BIC is an asymptotic approximation to a Bayesian approach to model selection.

BIC: Gaussian linear regression case

In the special case of a linear regression with Gaussian errors, the BIC simplifies to (ignoring constant terms)

$$\text{BIC} = \log \left(\frac{\text{RSS}}{N} \right) + \frac{\log(N)}{N} (p + 2)$$

<https://powcoder.com>

If we assume that the variance of the errors is known, we have instead

$$\text{BIC} = \frac{1}{\hat{\sigma}^2} \left(\frac{\text{RSS}}{N} + \frac{\log(N)}{N} \sigma^2 (p + 1) \right),$$

In this case the BIC is proportional to AIC and C_p , but with a $\log(N)$ penalty factor instead of 2.

Assignment Project Exam Help

Comparison of model selection

methods <https://powcoder.com>

Add WeChat powcoder

Model selection properties

Consistency. In a collection of models that includes the correct model, the probability that the model selection criterion chooses the correct one approaches one when $N \rightarrow \infty$.

Efficiency. The selected model predicts as well as the theoretically best model under consideration in terms of expected loss when $N \rightarrow \infty$.

Add WeChat powcoder

It is not possible to combine these properties (Claeskens and Hjort, 2008, Section 4.9).

Properties of model selection methods

Assignment Project Exam Help

LOOCV, AIC, and Mallows C_p Efficient but not consistent.

The efficiency follows because they construct unbiased estimators of the test error. However, they select models that are strictly more complex than the true model when $N \rightarrow \infty$.

BIC. Consistent under some conditions but not efficient. It often chooses models that are too simple because of its heavier penalty on complexity.

<https://powcoder.com>

Add WeChat powcoder

LOOCV, AIC and C_p

- LOOCV, AIC, and C_p are equivalent when $N \rightarrow \infty$. They will pick the same model in practice when N is large.

- In finite samples, we can view AIC and C_p as theoretical approximations to LOOCV.

- The advantage of AIC and C_p over LOOCV is mainly computational. CV should be preferred to AIC when the assumptions of the model (e.g., constant error variance) are likely to be wrong.

- LOOCV is universally applicable, while this is not the case for AIC and C_p .

Assignment Project Exam Help

Limitations of model selection
<https://powcoder.com>

Add WeChat powcoder

Limitations of model selection

- Standard statistical inference is no longer valid after model selection

Assignment Project Exam Help

- The reason is that standard inference assumes a fixed model, whereas model selection will by definition pick the specific model that best fits the sample. This will lead to optimistic estimates of sample variation based on the chosen model.

<https://powcoder.com>

Add WeChat powcoder

- In our context, the only way around this difficulty would be data splitting: using one part of the sample for model selection, and another for inference.

Assignment Project Exam Help

Model selection is an important tool in your data analysis process, but should not be a replacement to model building through EDA, diagnostics, and domain knowledge.

<https://powcoder.com>
Add WeChat powcoder

Assignment Project Exam Help

- What are the Akaike Information Criterion, Bayesian Information Criterion and Mallows's C_p ?
- What are the relationships between the above 3 metrics?
- Why is it incorrect to conduct statistical inference after model selection (using the same data)?

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

Optimism (optional)
<https://powcoder.com>

Add WeChat powcoder

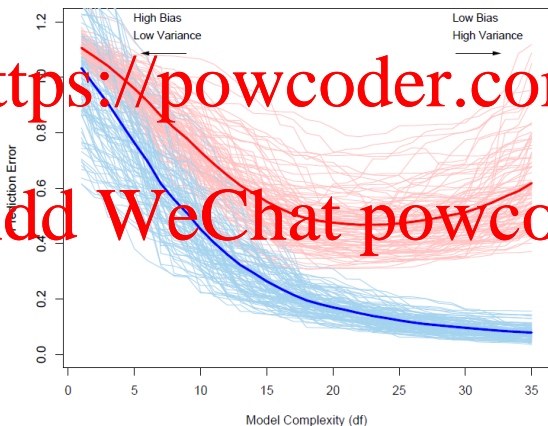
Optimism

Our objective in this section is to develop a better understanding of overfitting. This discussion will inform our understanding of analytical criteria in the next section.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



We define the **training error** as the empirical loss for the training

data,

Assignment Project Exam Help

$$\overline{\text{err}}_{\mathcal{D}} = \frac{1}{N} \sum_{i=1}^N L(Y_i, \hat{f}(\mathbf{x}_i)).$$

<https://powcoder.com>

We focus on our standard regression setting,

Add WeChat powcoder

$$\overline{\text{err}}_{\mathcal{D}} = \frac{1}{N} \sum_{i=1}^N (Y_i - \hat{f}(\mathbf{x}_i))^2,$$

which is RSS/N for linear regression estimated by least squares.

The expected prediction error (EPE) is

Assignment Project Exam Help

$$\begin{aligned}\text{EPE}(\mathbf{x}_0) &= E_{\mathcal{D}} \left[\left(Y_0 - \hat{f}(\mathbf{x}_0) \right)^2 \right] \\ &= E_{\mathcal{D}} \left[\left(f(\mathbf{x}_0) + \varepsilon_0 - \hat{f}(\mathbf{x}_0) \right)^2 \right]\end{aligned}$$

$$= \sigma^2 + E_{\mathcal{D}} \left[\left(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0) \right)^2 \right]$$

where \mathbf{x}_0 is fixed.

Now, consider the estimation error for a training case i .

$$\begin{aligned} E_{\mathcal{D}} \left(Y_i - \hat{f}(\mathbf{x}_i) \right)^2 &= E_{\mathcal{D}} \left[\left(f(\mathbf{x}_i) + \varepsilon_i - \hat{f}(\mathbf{x}_i) \right)^2 \right] \\ &= \sigma^2 + E_{\mathcal{D}} \left[\left(f(\mathbf{x}_i) - \hat{f}(\mathbf{x}_i) \right)^2 \right] - 2E_{\mathcal{D}} \left[\hat{f}(\mathbf{x}_i) \varepsilon_i \right] \\ &= \sigma^2 + E_{\mathcal{D}} \left[\left(f(\mathbf{x}_i) - \hat{f}(\mathbf{x}_i) \right)^2 \right] - 2\text{Cov}(\hat{f}(\mathbf{x}_i), \varepsilon_i) \end{aligned}$$

Unlike in the EPE, the last term appears because the estimator

$\hat{f}(\mathbf{x}_i)$ is a function of \mathcal{D} , which includes training case i itself.

Averaging over the data, the expected value of the training error is

Assignment Project Exam Help

$$E[\text{err}_{\mathcal{D}}] = \frac{1}{N} \sum_{i=1}^N E_{\mathcal{D}} \left[\left(Y_i - \hat{f}(\mathbf{x}_i) \right)^2 \right]$$

<https://powcoder.com>

Add WeChat powcoder

Because of the last term, the training error is not a good estimate of the expected prediction error.

Assignment Project Exam Help

We define the out of sample error as

$$\begin{aligned}\overline{\text{Err}}_{\text{out}} &= \frac{1}{N} \sum_{i=1}^N E \left[\left(Y_i^0 - \hat{f}(x_i) \right)^2 \right], \\ &= \sigma^2 + \frac{1}{N} \sum_{i=1}^N E_{\mathcal{D}} \left[\left(f(x_i) - \hat{f}(x_i) \right)^2 \right],\end{aligned}$$

where $Y_i^0 = f(x_i) + \epsilon_i^0$ indicates an independent case for a given x_i .

Add WeChat powcoder

Assignment Project Exam Help

The **optimism** of the training error is

$$\overline{L(\text{out})} - \overline{L(\text{Err}_D)} = \frac{2}{N} \sum_{i=1}^N \text{Cov}(\hat{f}(x_i), \varepsilon_i)$$

The more we overfit, the higher $\text{Cov}(\hat{f}(x_i), \varepsilon_i)$ will be, increasing the optimism.

Example: linear regression

For the linear regression model, we can show that

$$\text{Optimism} = \frac{2}{N} \sum_{i=1}^N \text{Cov}(\hat{f}(x_i), \varepsilon_i) = \frac{2}{N} \sigma^2 (p + 1)$$

Interpretation:

- The larger the sample size (N), the harder it is to overfit.
- The larger the variance of the errors (σ^2), the larger the overfitting.
- The optimism is proportional to the number of predictors.