

Basic Stats

Bill Last Updated:

26 January, 2020

Contents

Preface: Motivation	5
1 Logit and Probit	7
1.1 Logit	7
1.2 Probit	9
2 MLE	13
2.1 Basic idea of MLE	13
2.2 Coin flip example, probit, and logit	14
2.3 Further on logit	16
2.4 References	17
3 Score, Gradient and Jacobian	19
3.1 Score	19
3.2 Fisher scoring	20
3.3 Gradient and Jacobian	20
3.4 Hessian and Fisher Information	21
4 Canonical link function	23
5 Ordinary Least Squares (OLS)	25
5.1 Taylor series	27
5.2 References	27

6	Measure Theory	29
6.1	Part 1	29
6.2	Part 2	30
6.3	Part 3	31
6.4	Part 4	32
6.5	Part 5	32
6.6	Part 6	33

Preface: Motivation

All the notes I have done here are about basic stats. While I have tried my best, probably there are still some typos and errors. Please feel free to let me know in case you find one. Thank you!

Chapter 1

Logit and Probit

1.1 Logit

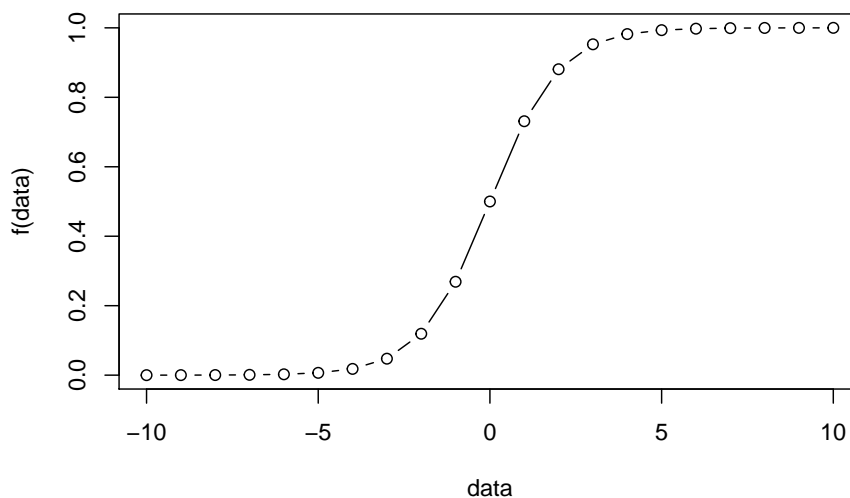
$$f(x) = \log\left(\frac{p(y=1)}{1-p(y=1)}\right)$$

The basic idea of logistic regression:

$$p(y=1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

Thus, $e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}$ can be from $-\infty$ to $+\infty$, and $p(y=1)$ will be always within the range of $(0, 1)$.

```
f<-function(x){exp(x)/(1+exp(x))}  
data<-seq(-10,10,1)  
plot(data,f(data),type = "b")
```



We can also write the function into another format as follows:

$$\log \frac{p(y=1)}{1-p(y=1)} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

Thus, we know that the regression coefficients of β_i actually change the “log-odds” of the event. Of course, note that the magnitude of β_i is dependent upon the units of x_i .

The following is an example testing whether that home teams are more likely to win in NFL games. The results show that the odd of winning is the same for both home and away teams.

```
mydata = read.csv(url('https://raw.githubusercontent.com/nfl-football-ops/Big-Data-Bow
mydata$result_new<-ifelse(mydata$HomeScore>mydata$VisitorScore,1,0)
summary(mydata$result_new)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.0000 0.0000 0.0000 0.4945 1.0000 1.0000
```

```
mylogit1 = glm(result_new~1, family=binomial, data=mydata)
summary(mylogit1)
```

```
##
## Call:
```



```
## glm(formula = result_new ~ 1, family = binomial, data = mydata)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.168  -1.168  -1.168   1.187   1.187
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.02198    0.20967  -0.105   0.917
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 126.14  on 90  degrees of freedom
## Residual deviance: 126.14  on 90  degrees of freedom
## AIC: 128.14
##
## Number of Fisher Scoring iterations: 3
```

1.2 Probit

As noted above, logit $f(x) = \log\left(\frac{p(y=1)}{1-p(y=1)}\right)$ provides the resulting range of $(0, 1)$. Another way to provide the same range is through the cdf of normal distribution. The following R code is used to illustrate this process.

```
data2<-seq(-5,5,1)
plot(data2,pnorm(data2),type = "b")
```



Thus, the cdf of normal distribution can be used to indicate the probability of $p(y = 1)$.

$$\Phi(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n) = p(y = 1)$$

Similar to logit model, we can also write the inverse function of the cdf to get the function that can be from $-\infty$ to $+\infty$.

$$\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n = \Phi^{-1}(p(y = 1))$$

Thus, for example, if $X\beta = -2$, based on $\Phi(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n) = p(y = 1)$ we can get that the $p(y = 1) = 0.023$.

In contrast, if $X\beta = 3$, the $p(y = 1) = 0.999$.

```
pnorm(-2)
```

```
## [1] 0.02275013
```

```
pnorm(3)
```

```
## [1] 0.9986501
```

Let's assume that there is a latent variable called Y^* such that

$$Y^* = X\beta + \epsilon, \epsilon \sim N(0, \sigma^2)$$

You could think of Y^* as a kind of “proxy” between $X\beta + \epsilon$ and the observed $Y(1 \text{ or } 0)$. Thus, we can get the following. Note that, it does not have to be zero, and can be any constant.

$$Y^* = \begin{cases} 0 & \text{if } y_i^* \leq 0 \\ 1 & \text{if } y_i^* > 0 \end{cases}$$

Thus,

$$y_i^* > 0 \Rightarrow \beta' X_i + \epsilon_i > 0 \Rightarrow \epsilon_i > -\beta' X_i$$

Thus, we can write it as follows. Note that $\frac{\epsilon_i}{\sigma} \sim N(0, 1)$

$$p(y = 1|x_i) = p(y_i^* > 0|x_i) = p(\epsilon_i > -\beta' X_i) = p\left(\frac{\epsilon_i}{\sigma} > \frac{-\beta' X_i}{\sigma}\right) = \Phi\left(\frac{\beta' X_i}{\sigma}\right)$$

We thus can get:

$$p(y = 0|x_i) = 1 - \Phi\left(\frac{\beta' X_i}{\sigma}\right)$$

For $p(y = 1|x_i) = \Phi\left(\frac{\beta' X_i}{\sigma}\right)$, we can not really estimate both β and σ as they are in a ratio. We can assume $\sigma = 1$, then $\epsilon \sim N(0, 1)$. We know y_i and x_i since we observe them. Thus, we can write it as follows.

$$p(y = 1|x_i) = \Phi(\beta' X_i)$$

Chapter 2

MLE

2.1 Basic idea of MLE

Suppose that we flip a coin, $y_i = 0$ for tails and $y_i = 1$ for heads. If we get p heads from n trials, we can get the proportion of heads is p/n , which is the sample mean. If we do not do any further calculation, this is our best guess.

Suppose that the true probability is ρ , then we can get:

$$\mathbf{L}(y_i) = \begin{cases} \rho & y_i = 1 \\ 1 - \rho & y_i = 0 \end{cases}$$

Thus, we can also write it as follows.

$$\mathbf{L}(y_i) = \rho^{y_i} (1 - \rho)^{1-y_i}$$

Thus, we can get:

$$\prod \mathbf{L}(y_i|\rho) = \rho^{\sum y_i} (1 - \rho)^{\sum (1-y_i)}$$

Further, we can get a log-transformed format.

$$\log(\prod \mathbf{L}(y_i|\rho)) = \sum y_i \log \rho + \sum (1 - y_i) \log(1 - \rho)$$

To maximize the log-function above, we can calculate the derivative with respect to ρ .

$$\frac{\partial \log(\prod \mathbf{L}(y_i|\rho))}{\partial \rho} = \sum y_i \frac{1}{\rho} - \sum (1 - y_i) \frac{1}{1 - \rho}$$

Set the derivative to zero and solve for ρ , we can get

$$\begin{aligned}
& \sum y_i \frac{1}{\rho} - \sum (1 - y_i) \frac{1}{1 - \rho} = 0 \\
& \Rightarrow (1 - \rho) \sum y_i - \rho \sum (1 - y_i) = 0 \\
& \Rightarrow \sum y_i - \rho \sum y_i - n\rho + \rho \sum y_i = 0 \\
& \Rightarrow \sum y_i - n\rho = 0 \\
& \Rightarrow \rho = \frac{\sum y_i}{n} = \frac{p}{n}
\end{aligned}$$

Thus, we can see that the ρ maximizing the likelihood function is equal to the sample mean.

2.2 Coin flip example, probit, and logit

In the example above, we are not really trying to estimate a lot of regression coefficients. What we are doing actually is to calculate the sample mean, or intercept in the regression sense. What does it mean? Let's use some data to explain it.

Suppose that we flip a coin 20 times and observe 8 heads. We can use the R's `glm` function to estimate the ρ . If the result is consistent with what we did above, we should observe that the *cdf* of the estimate of β_0 (i.e., intercept) should be equal to $8/20 = 0.4$.

```
coins<-c(rep(1,times=8),rep(0,times=12))
table(coins)
```

```
## coins
##  0  1
## 12  8
```

```
coins<-as.data.frame(coins)
```

2.2.1 Probit

```
probitresults <- glm(coins ~ 1, family = binomial(link = "probit"), data = coins)
probitresults
```

```
##
## Call:  glm(formula = coins ~ 1, family = binomial(link = "probit"),
##       data = coins)
##
## Coefficients:
## (Intercept)
##      -0.2533
##
## Degrees of Freedom: 19 Total (i.e. Null);  19 Residual
## Null Deviance:      26.92
## Residual Deviance: 26.92    AIC: 28.92
```

```
pnorm(probitresults$coefficients)
```

```
## (Intercept)
##           0.4
```

As we can see the intercept is -0.2533 , and thus $\Phi(-0.2533471) = 0.4$

2.2.2 Logit

We can also use logit link to calculate the intercept as well. Recall that

$$p(y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

Thus,

$$p(y = 1) = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$$

```
logitresults <- glm(coins ~ 1, family = binomial(link = "logit"), data = coins)
logitresults$coefficients
```

```
## (Intercept)
##      -0.4054651
```

```
exp(logitresults$coefficients)/(1+exp(logitresults$coefficients))
```

```
## (Intercept)
##           0.4
```

Note that, the default link for the binomial in the glm function is logit.

2.3 Further on logit

The probability of $y = 1$ is as follows:

$$p = p(y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

Thus, the likelihood function is as follows:

$$\begin{aligned} L &= \prod p^{y_i} (1-p)^{1-y_i} = \prod \left(\frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} \right)^{y_i} \left(\frac{1}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}} \right)^{1-y_i} \\ &= \prod (1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)})^{-y_i} (1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n})^{-(1-y_i)} \end{aligned}$$

Thus, the log-likelihood is as follows:

$$\log L = \sum (-y_i \cdot \log(1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}) - (1 - y_i) \cdot \log(1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}))$$

Typically, optimisers minimize a function, so we use negative log-likelihood as minimising that is equivalent to maximising the log-likelihood or the likelihood itself.

#Source of R code: <https://www.r-bloggers.com/logistic-regression/>

```
mle.logreg = function(fmla, data)
{
  # Define the negative log likelihood function
  logl <- function(theta,x,y){
    y <- y
    x <- as.matrix(x)
    beta <- theta[1:ncol(x)]

    # Use the log-likelihood of the Bernoulli distribution, where p is
    # defined as the logistic transformation of a linear combination
    # of predictors, according to logit(p)=(x%*%beta)
    loglik <- sum(-y*log(1 + exp(-(x%*%beta))) - (1-y)*log(1 + exp(x%*%beta)))
    return(-loglik)
  }

  # Prepare the data
  outcome = rownames(attr(terms(fmla),"factors"))[1]
  dfrTmp = model.frame(data)
  x = as.matrix(model.matrix(fmla, data=dfrTmp))
}
```



```

y = as.numeric(as.matrix(data[,match(outcome,colnames(data))]))

# Define initial values for the parameters
theta.start = rep(0,(dim(x)[2]))
names(theta.start) = colnames(x)

# Calculate the maximum likelihood
mle = optim(theta.start,logl,x=x,y=y, method = 'BFGS', hessian=T)
out = list(beta=mle$par,vcov=solve(mle$hessian),ll=2*mle$value)
}

mydata = read.csv(url('https://stats.idre.ucla.edu/stat/data/binary.csv'))
mylogit1 = glm(admit~gre+gpa+as.factor(rank), family=binomial, data=mydata)

mydata$rank = factor(mydata$rank) #Treat rank as a categorical variable
fmla = as.formula("admit~gre+gpa+rank") #Create model formula
mylogit2 = mle.logreg(fmla, mydata) #Estimate coefficients

print(cbind(coef(mylogit1), mylogit2$beta))

##                [,1]      [,2]
## (Intercept)    -3.989979073 -3.772676422
## gre            0.002264426  0.001375522
## gpa            0.804037549  0.898201239
## as.factor(rank)2 -0.675442928 -0.675543009
## as.factor(rank)3 -1.340203916 -1.356554831
## as.factor(rank)4 -1.551463677 -1.563396035

```

2.4 References

http://www.columbia.edu/~so33/SusDev/Lecture_9.pdf

Chapter 3

Score, Gradient and Jacobian

3.1 Score

The score is the gradient (the vector of partial derivatives) of $\log L(\theta)$, with respect to an m -dimensional parameter vector θ .

$$S(\theta) = \frac{\partial \ell}{\partial \theta}$$

Typically, they use ∇ to denote the partial derivative.

$$\nabla \ell$$

Such differentiation will generate a $m \times 1$ row vector, which indicates the sensitivity of the likelihood.

Quote from Steffen Lauritzen's slides: "Generally the solution to this equation must be calculated by iterative methods. One of the most common methods is the Newton–Raphson method and this is based on successive approximations to the solution, using Taylor's theorem to approximate the equation."

For instance, using logit link, we can get the first derivative of log likelihood logistic regression as follows. We can not really find β easily to make the equation to be 0.

$$\begin{aligned}\frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^n x_i^T \left[y_i - \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \right] \\ &= \sum_{i=1}^n x_i^T [y_i - \hat{y}_i]\end{aligned}$$

3.2 Fisher scoring

[I will come back to this later.]

<https://www2.stat.duke.edu/courses/Fall00/sta216/handouts/diagnostics.pdf>

<https://stats.stackexchange.com/questions/176351/implement-fisher-scoring-for-linear-regression>

3.3 Gradient and Jacobian

Remarks: This part discusses gradient in a more general sense.

When $f(x)$ is only in a single dimension space:

$$\mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla f(x) = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

When $f(x)$ is only in a m-dimension space (i.e., Jacobian): $\mathbb{R}^n \rightarrow \mathbb{R}^m$

$$Jac(f) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \frac{\partial f_m}{\partial x_3} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

For instance,

$$\mathbb{R}^n \rightarrow \mathbb{R}:$$

$$f(x, y) = x^2 + 2y$$

$$\nabla f(x, y) = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right] = [2x, 2]$$

$$\mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f(x, y) = (x^2 + 2y, x^3)$$

$$Jac(f) = \begin{bmatrix} 2x & 2 \\ 2x^2 & 0 \end{bmatrix}$$

3.4 Hessian and Fisher Information

Hessian matrix or Hessian is a square matrix of second-order partial derivatives of a scalar-valued function, or scalar field.

$\mathbb{R}^n \rightarrow \mathbb{R}$

$$Hessian = \nabla^2(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_1 \partial x_3} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \frac{\partial^2 f}{\partial x_2 \partial x_3} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \frac{\partial^2 f}{\partial x_3 \partial x_1} & \frac{\partial^2 f}{\partial x_3 \partial x_2} & \frac{\partial^2 f}{\partial x_3^2} & \cdots & \frac{\partial^2 f}{\partial x_3 \partial x_n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \frac{\partial^2 f}{\partial x_n \partial x_3} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

As a special case, in the context of logit:

Suppose that the log likelihood function is $\ell(\theta)$. θ is a m dimension vector.

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \cdots \\ \theta_m \end{bmatrix}$$

$$Hessian = \nabla^2(\ell) = \begin{bmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_3} & \cdots & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_m} \\ \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_2^2} & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_3} & \cdots & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_m} \\ \frac{\partial^2 \ell}{\partial \theta_3 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_3 \partial \theta_2} & \frac{\partial^2 \ell}{\partial \theta_3^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_3 \partial \theta_m} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 \ell}{\partial \theta_m \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_m \partial \theta_2} & \frac{\partial^2 \ell}{\partial \theta_m \partial \theta_3} & \cdots & \frac{\partial^2 \ell}{\partial \theta_m \partial \theta_m} \end{bmatrix}$$

“In statistics, the observed information, or observed Fisher information, is the negative of the second derivative (the Hessian matrix) of the “log-likelihood” (the logarithm of the likelihood function). It is a sample-based version of the Fisher information.” (Direct quote from Wikipedia.)

Thus, the observed information matrix:

$$-Hessian = -\nabla^2(\ell)$$

Expected (Fisher) information matrix:

$$E[-\nabla^2(\ell)]$$

Chapter 4

Canonical link function

Inspired by a Stack Exchange post, I created the following figure:

$$\frac{\text{Parameter}}{\theta} \rightarrow \gamma'(\theta) = \mu \rightarrow \frac{\text{Mean}}{\mu} \rightarrow g(\mu) = \eta \rightarrow \frac{\text{Linear predictor}}{\eta}$$

For the case of n time Bernoulli (i.e., Binomial), its canonical link function is logit. Specifically,

$$\frac{\text{Parameter}}{\theta = \beta^T x_i} \rightarrow \gamma'(\theta) = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \rightarrow \frac{\text{Mean}}{\mu = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}} \rightarrow g(\mu) = \log \frac{\frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}}{1 - \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}} \rightarrow \frac{\text{Linear predictor}}{\eta = \beta^T x_i}$$

Thus, we can see that,

$$\theta \equiv \eta$$

The link function $g(\mu)$ relates the linear predictor $\eta = \beta^T x_i$ to the mean μ .

Remarks:

- (1) Parameter is $\theta = \beta^T x_i$ (Not μ !).
- (2) $\mu = p(y = 1) = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$ (Not logit!).
- (3) Link function (i.e., $g(\mu)$) = logit = logarithm of odds = $\log \frac{\text{Event-Happened}}{\text{Event-Not-Happened}}$.
- (4) $g(\mu) = \log \frac{\mu}{1-\mu} = \beta^T x_i$. Thus, link function = linear predictor = log odds!

- (5) Quote from the Stack Exchange post “Newton Method and Fisher scoring for finding the ML estimator coincide, these links simplify the derivation of the MLE.”

(Recall, we know that μ or $p(y = 1)$ is the mean function. Recall that, n trails of coin flips, and get p heads. Thus $\mu = \frac{p}{n}$.)

Chapter 5

Ordinary Least Squares (OLS)

Suppose we have n observation, and m variables.

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2m} \\ \dots & & & & \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix}$$

Thus, we can write it as the following n equations.

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_m x_{1m}$$

$$y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m}$$

$$y_3 = \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} + \dots + \beta_m x_{3m}$$

...

$$y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm}$$

We can combine all the n equations as the following one:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_m x_{im} (i \in [1, n])$$

We can further rewrite it as a matrix format as follows.

$$y = X\beta$$

Where,

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ \dots \\ y_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2m} \\ \dots & & & & & \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \dots \\ \beta_m \end{bmatrix}$$

Since later we need the inverse of X , we need to make it into a square matrix.

$$X^T y = X^T X \hat{\beta} \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

We can use R to implement this calculation. As we can see, there is no need to do any iterations at all, but rather just pure matrix calculation.

```
X<-matrix(rnorm(1000),ncol=2) # we define a 2 column matrix, with 500 rows
X<-cbind(1,X) # add a 1 constant
beta_true<-c(2,1,2) # True regression coefficients
beta_true<-as.matrix(beta_true)
y=X%%beta_true+rnorm(500)
```

```
transposed_X<-t(X)
beta_hat<-solve(transposed_X%%X)%%transposed_X%%y
beta_hat
```

```
##           [,1]
## [1,] 2.0081310
## [2,] 0.9689952
## [3,] 2.0324824
```

Side Notes The function of `as.matrix` will automatically make `c(2,1,2)` become the dimension of 3×1 , you do not need to transpose the β .

5.1 Taylor series

$$\begin{aligned} f(x)|_a &= f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f'(a)}{2!}(x-a)^2 + \frac{f''(a)}{3!}(x-a)^3 + \dots \\ &= \sum_{n=0}^{\infty} \frac{f^n(a)}{n!}(x-a)^n \end{aligned}$$

For example:

$$\begin{aligned} e^x|_{a=0} &= e^a + \frac{e^a}{1!}(x-a) + \frac{e^a}{2!}(x-a)^2 + \dots + \frac{e^a}{n!}(x-a)^n \\ &= 1 + \frac{1}{1!}x + \frac{1}{2!}x^2 + \dots + \frac{1}{n!}x^n \end{aligned}$$

if $x = 2$

$$e^2 = 7.389056$$

$$e^2 \approx 1 + \frac{1}{1!}x = 1 + \frac{1}{1!}2 = 3$$

$$e^2 \approx 1 + \frac{1}{1!}x + \frac{1}{2!}x^2 = 1 + \frac{1}{1!}2 + \frac{1}{2!}2 = 5 \dots$$

$$e^2 \approx 1 + \frac{1}{1!}x + \frac{1}{2!}x^2 + \frac{1}{3!}x^2 + \frac{1}{4!}x^2 + \frac{1}{5!}x^2 = 7.2666\dots$$

5.2 References

1. Steffen Lauritzen's slides:

<http://www.stats.ox.ac.uk/~steffen/teaching/bs2HT9/scoring.pdf>

2. The Stack Exchange post:

<https://stats.stackexchange.com/questions/40876/what-is-the-difference-between-a-link-function-and-a-canonical-link-function>

3. Wikipedia for OLS

https://en.wikipedia.org/wiki/Ordinary_least_squares

4. Gradient and Jacobian

<https://math.stackexchange.com/questions/1519367/difference-between-gradient-and-jacobian>

https://www.youtube.com/watch?v=3xVMVT-2_t4

<https://math.stackexchange.com/questions/661195/what-is-the-difference-between-the-gradient-and-the-directional-derivative>

5. Hessian

https://en.wikipedia.org/wiki/Hessian_matrix

6. Observed information

https://en.wikipedia.org/wiki/Observed_information

7. Fisher information

https://people.missouristate.edu/songfengzheng/Teaching/MTH541/Lecture%20notes/Fisher_info.pdf

8. Link function

https://en.wikipedia.org/wiki/Generalized_linear_model#Link_function

<https://stats.stackexchange.com/questions/40876/what-is-the-difference-between-a-link-function-and-a-canonical-link-function>

Chapter 6

Measure Theory

Instead of density, give probability to every event.

Definition

Let Ω be set [state space]= set of all outcomes.

Each $\omega \in \Omega$ is outcome.

Each $A \in \mathcal{A}$ is event.

For rolling a die, $\Omega = 1, 2, \dots, 6$, whereas $A = 1, 3, 6$. Thus, $P(A) = \frac{1}{2}$.

$P : 2^\Omega \rightarrow [0, 1]$

$$(1) \ P(\Omega) = 1, P(\emptyset) = 0$$

$$(2) \ P(A) \in [0, 1]$$

$$(3) \ \text{If } A \cap B = \emptyset, \text{ then } P(A \cup B) = P(A) + P(B)$$

$$(4) \ \text{If } A_i \text{ are mutually disjoint, then } P(\cup A_i) = \sum P(A_i)$$

6.1 Part 1

$P(X)$: Power set X

Example: the set is $X = \{a, b\}$, and thus the power set is $p(x) = \{\emptyset, X, a, b\}$

Definition: $\mathcal{A} \subseteq P(X)$ is called a σ -algebra.

- (a) σ -algebra at least has two elements, namely empty set and full set: $\emptyset, X \in \mathcal{A}$

$$(b) A \in \mathcal{A} \rightarrow A^c \in \mathcal{A}$$

$$(C) A_i \in \mathcal{A} \rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$$

____ $A \in \mathcal{A}$ is called an \mathcal{A} -measurable set.

Example of σ -algebra:

$$(1) \mathcal{A} = \{\emptyset, X\}$$

$$(2) \mathcal{A} = P(X)$$

6.2 Part 2

Easy to show:

\mathcal{A}_i σ -algebra on X , in $I(indexset)$. Then, $\bigcap \mathcal{A}_i$ is also a σ -algebra on X .

Definition: For $M \subseteq P(X)$, there is a smallest σ -algebra that contains M . That is, σ -algebra is generated by M .

$$\sigma(M) := \bigcap_{\mathcal{A} \supseteq M} \mathcal{A}$$

Example:

$$X = \{a, b, c, d\}, M = \{\{a\}, \{b\}\}$$

Note that, M is not a σ -algebra yet, and the smallest σ -algebra containing M :

$$\sigma(M) = \{\emptyset, X = \{a, b, c, d\}, \{a\}, \{b\}, \{a, b\}, \{b, c, d\}, \{a, c, d\}, \{c, d\}\}$$

Definition:

Let X be a topological space

(Let X be a metric space, or let X be a subset of \mathbb{R}^n)

$B(X)$ Borel σ -algebra on X : The σ -algebra generated by open sets.

$$B(X) := \sigma(\mathfrak{I})$$

6.3 Part 3

What is a measure?

A measure is a σ -algebra. It can be either a power set or a subset.

Definition:

(X, \mathcal{A}) measurable space. \mathcal{A} is the subset of the set X .

A map $\mu : \mathcal{A} \rightarrow [0, \infty]$ is called a measure if it satisfies:

- (a) $\mu(\emptyset) = 0$
- (b) $\mu(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \mu(A_i)$ with $A_i \cap A_j = \emptyset$ for all $A_i \in \mathcal{A}$

Definition: Measure space

$$(X, \mathcal{A}, \mu)$$

Examples:

$$X, \mathcal{A} = P(X)$$

- (a) For $A \in \mathcal{A}$, counting measure:

$$\begin{cases} \text{Count}(A) & \text{if } A \text{ has finite elements} \\ \infty & \text{else} \end{cases}$$

Calculation rules in $[0, \infty]$:

- (1) $X + \infty := \infty$ for $X \in [0, \infty]$
- (2) $X \cdot \infty := \infty$ for $X \in (0, \infty]$
- (3) $0 \cdot \infty := 0$ (for $X \in (0, \infty]$) in most cases in measure theory!

- (b) Dirac measure for $p \in X$

$$\delta_p(A) = \begin{cases} 1 & , p \in A \\ 0 & , \text{else} \end{cases}$$

- (c) We search a measure on $X = \mathbb{R}^n$

- (1) $\mu([0, 1]^n) = 1$
- (2) $\mu(x + A) = \mu(A)$ for all $x \in \mathbb{R}^n$

6.4 Part 4

Search measure μ on $p(\mathbb{R})$ with

- (1) $\mu([a, b]) = b - a, b > a$
- (2) $\mu(X + A) = \mu(A), A \in p(\mathbb{R}), X \in \mathbb{R}$

$\rightarrow \mu$ does not exist.

Claim: Let μ be a measure on $p(\mathbb{R})$ with $\mu([0, 1]) < \infty$ and (2) $\rightarrow \mu = 0$

6.5 Part 5

Definition:

$(\Omega_1, \mathcal{A}_1), (\Omega_2, \mathcal{A}_2)$ measurable spaces

$f : \Omega_1 \rightarrow \Omega_2$ measurable (*w.r.t* $\mathcal{A}_1, \mathcal{A}_2$) if $f^{-1}(\mathcal{A}_2) \in \mathcal{A}_1$ for all $\mathcal{A}_2 \in \mathcal{A}_2$.

Examples:

- (1) $(\Omega, \mathcal{A}), (\mathbb{R}, B(\mathbb{R}))$: characteristic function (also indicator function)

$$\chi_A : \Omega \rightarrow \mathbb{R}$$

$$\chi_A(\omega) := \begin{cases} 1 & , \omega \in A \\ 0 & , else \end{cases}$$

For all measurable $A \in \mathcal{A}$, χ_A is a measurable map.

$$\begin{aligned} \chi_A(\emptyset) &= \emptyset, \chi_A^{-1}(\mathbb{R}) = \Omega \\ \chi_A(\{1\}) &= A, \chi_A^{-1}(\mathbb{R}^c) = A^c \end{aligned}$$

- (2) $(\Omega_1, \mathcal{A}_1), (\Omega_2, \mathcal{A}_2), (\Omega_3, \mathcal{A}_3)$ measurable spaces

$$\Omega_1 \rightarrow \Omega_2 \rightarrow \Omega_3$$

If both $\Omega_1 \rightarrow \Omega_2$ and $\Omega_2 \rightarrow \Omega_3$ are measurable, then $\Omega_1 \rightarrow \Omega_3$ is measurable as well.

6.6 Part 6

Lebesgue integral

Lebesgue integral for step functions.

$$\text{Measure space}(X, \mathcal{A}, \mu)$$

Where,

- (1) X is a set.
- (2) \mathcal{A} : collection of subset of X : σ -algebra.
- (3) $\mu: \mathcal{A} \rightarrow [0, \infty]$.

In order to integrate some special function, we need measurable maps:

$$\text{Measure maps } f : X \rightarrow \mathbb{R}$$

Where, there is σ -algebra \mathcal{A} on the X , whereas σ -algebra \mathcal{B} for \mathbb{R} . In other words, we can get:

$$f^{-1}(E) \in \mathcal{A} \text{ for all Borel sets } E \subseteq \mathbb{R}$$

For example:

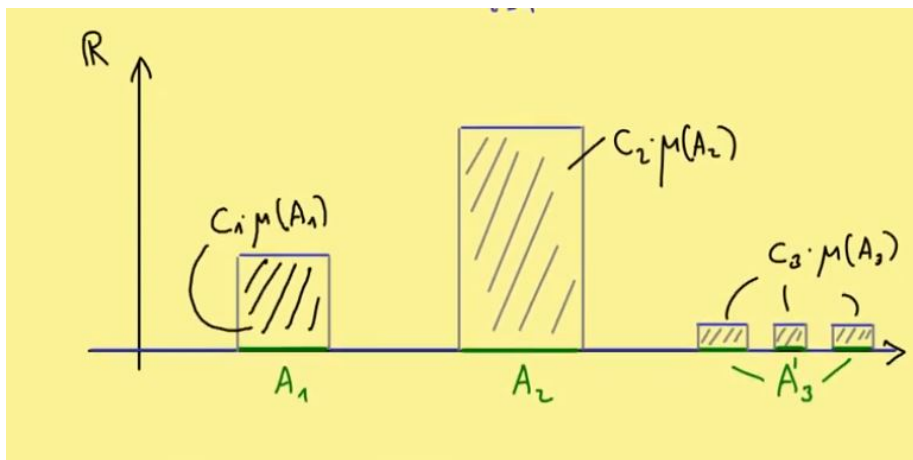
$$\chi_A : X \rightarrow \mathbb{R}, A \in \mathcal{A}$$

$$I(\chi_A) := \mu(A)$$

Simple functions (step functions, staircase functions, ...)

For $A_1, \dots, A_N \in \mathcal{A}$ and $c_1, \dots, c_n \in \mathbb{R}$

$$f(x) = \sum_{i=1}^n c_i \cdot \chi_{A_i}(x)$$

**Definition**

$$S^+ := \{f : X \rightarrow \mathbb{R} \mid f : \text{simple function}, f \geq 0\}$$

Properties:

- (1) $I(\alpha f + \beta g) = \alpha I(f) + \beta I(g)$, $\alpha, \beta \geq 0$
- (2) $f \leq g \rightarrow I(f) \leq I(g)$ (monotonicity)

Definitions

$f : X \rightarrow [0, \infty)$ measurable

$\int f d(\mu) := \sup\{I(h) \mid h \in S^+, h \leq f\}$, in which, $h = \sum_{i=1}^n c_i x_{A_i}$. We call this Lebesgue integral of f w.r.t μ . f is called μ -integral if $\int_x f d\mu < \infty$

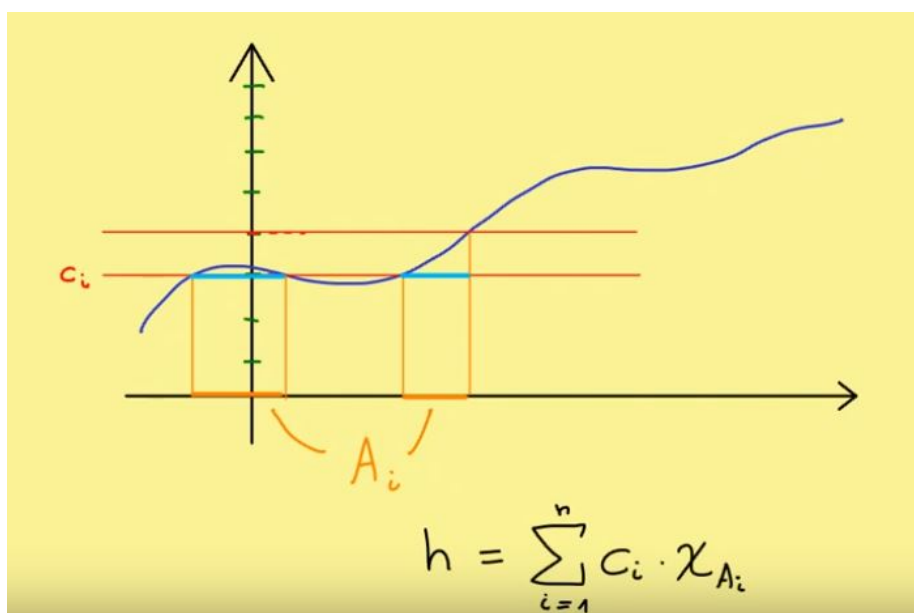


Figure 6.1: