



Model Assessment and Selection

By Mehrab Atighi

مقدمه‌ای بر یادگیری آماری

در این درس می‌خواهیم به این دو مسئله در مسایل رگرسیون بپردازیم

Model selection

- estimating the performance **of different models in** order to choose the **best one**.

Model assessment

- having chosen a final model, **estimating its prediction error** (generalization error) on **new data**

انواع یادگیری آماری

Model selection

Train

Validation

Test

- It is difficult to give a general rule on how to choose the number of observations in each of the three parts, as this depends on the signal to noise ratio in the data and the training sample size. A typical split might be 50% for training, and 25% each for validation and testing

Model assessment

- The methods of this chapter approximate the validation step either analytically (**AIC**, **BIC**, **MDL**, **SRM**) or by efficient sample re-use (crossvalidation and the bootstrap).

داده ها و بخش بندی کردن آنها

داده های آموزشی: این داده های supervise هستند و ما مدل خود را براساس این داده ها تقسیم بندی می کنیم.

داده های آزمایشی: این داده ها نیز برای تست کردن و بدست آوردن دقت مدل ما یا همان خطای مدل رگرسیونی ما مورد استفاده قرار می گیرند.

باید به این نیز توجه داشته باشیم که داده های آموزشی و داده های آزمایشی ما شامل زوج های p متغیره از داده ها و متغیر پاسخ و متغیرهای مستقل ما نیز هستند که بطور کلی داده هایمان به این دو تقسیم می شوند.



```
sample <- sample(c(TRUE,FALSE),nrow(Data),replace = T,prob = c(0.6,0.4))
train <- Data[sample,]
test <- Data[!sample,]
```

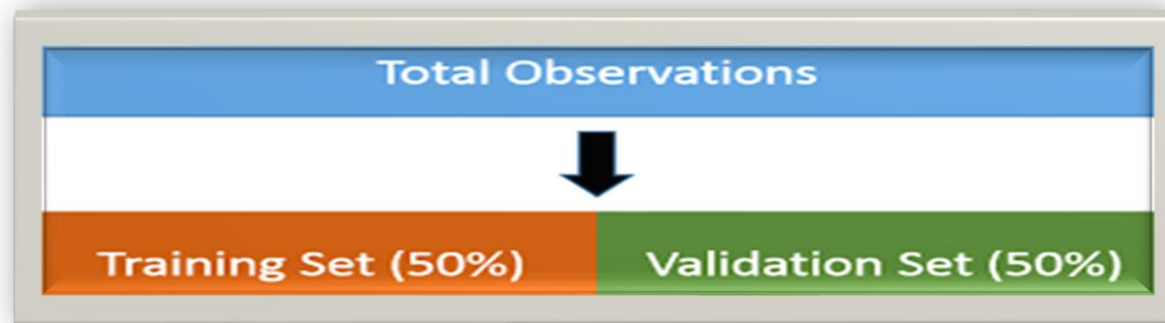
Validation set Approach

ما در این روش زمانی که داده‌هایمان را به داده‌های آموزش و آزمون تقسیم بندی می‌کنیم بصورت ۶۰-۴۰ یا ۷۰-۳۰ این کار را انجام می‌دهیم یعنی اینکه مثلاً ۶۰ درصد داده‌های ما آموزش باشند، که به کمک آنها مدل را بسازیم و مدل ساخته شده را روی ۴۰ درصد مابقی که شامل داده‌های تست یا آزمون ما هستند، تست کرده و میزان خطای آنرا براساس معیار موردنظر بررسی بکنیم.

```
sample <- sample(c(TRUE, FALSE), nrow(Data), replace = T, prob = c(0.6, 0.4))  
train <- Data[sample, ]  
test <- Data[!sample, ]
```

برای داشتن یک مدل رگرسیون درجه دو سه و یا بیشتر بصورت روبرو انجام می‌دهیم:

```
lm(y~poly(x, deg))
```



Cross validation

• ارزیابی متقابل:

چرا روش‌های باز نمونه‌گیری را بکار ببریم؟

گاهی اوقات که ما می‌آییم و بصورت تصادفی و مکرر از داده‌های آموزشی خود مجموعه‌های مختلفی را انتخاب می‌کنیم و هربار اطلاعات جدیدتری به ما ارائه می‌شود. و ما دنبال دو هدف ارزیابی و انتخاب مدل خود هستیم.

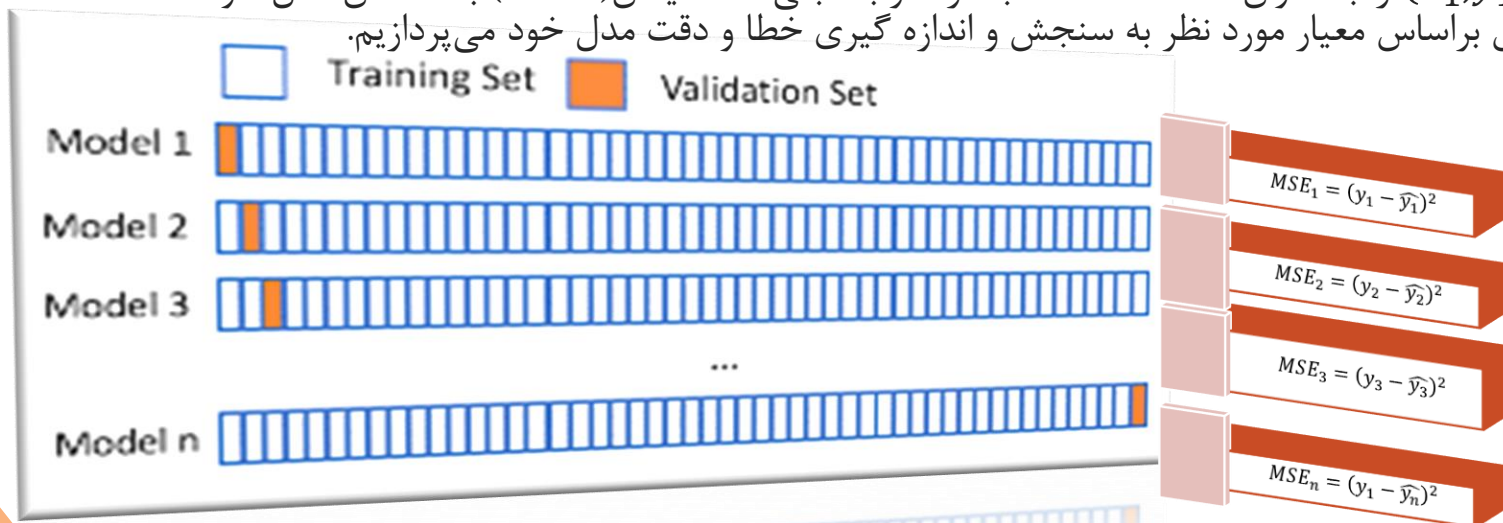
- Leave one out cross validation
- K fold cross validation
- Cross validation for classification problem miss rate

Leave one out cross validation

فرض کنید که تعداد داده‌های ما خیلی زیاد نیست ما از این روش استفاده می‌کنیم با فرض داشتن n داده داریم:



در این روش ما (x_1, y_1) را به عنوان داده تست انتخاب کرده و با مابقی داده‌هایمان ($n-1$ تا) به ساختن مدل خود می‌پردازیم و سپس براساس معیار مورد نظر به سنجش و اندازه‌گیری خطا و دقت مدل خود می‌پردازیم.



حال این n تا معیار نیز می‌توانند به ازای هر درجه مختلفی از مدل رگرسیون مثلا درجه یک دو سه و ... تا n داده باشند و داریم:

$$Cv_n = \frac{1}{n} \sum_{i=1}^n MSE_i$$

Leave one out cross validation

- #R codes:

```
glm.fit <- glm(y~x, data = ..., ...)
```

```
loocv.err <- cv.glm(data, glm.fit)
```

```
str(loocv.err)
```

```
loocv.err$ delta[n]
```

همان mse_n ما می باشد:

در این روش MSE ما بسیار کم می شود ولی از نظر زمانی هرچه تعداد داده های ما بیشتر شود، زمان بر می شود و محاسبات زیادتر و سنگین تر خواهد داشت.

ارزیابی متقابل k تایی:

The diagram illustrates the K-fold cross-validation process. It shows a large bar at the top labeled "All Data". Below this, the data is divided into "Training Data" (blue bar) and "Test Data" (orange bar). The "Training Data" is further divided into five "FOLD" sections (FOLD 1 to FOLD 5). The "Test Data" is also divided into five "FOLD" sections (FOLD 1 to FOLD 5). The process is repeated for each fold, with the training data being split into K-1 folds and the test data being the remaining 1 fold. This results in K split scores (Split 1 score to Split 5 score). The final score is the average of all split scores.

All Data					
Training Data					Test Data
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	Split 1 score
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	Split 2 score
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	Split 3 score
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	Split 4 score
FOLD 1	FOLD 2	FOLD 3	FOLD 4	FOLD 5	Split 5 score
Final score					Average of all split scores

K Fold Cross Validation

#R codes for K Fold Cross Validation:

در کدنویسی همان تابع *cv.glm* را استفاده می‌کنیم با این تفاوت که در اینجا باید میزان *k* را ذکر کنیم:

```
kv.fit <- -lm(y~x, data = ... )  
cv.glm(data,kv.fit)$delta[n]      =  $MSE_n$   
kf cv - err <- -function(x){  
  fit <- -lm(y~x, data = ... )  
  cv.glm(data , fit, k = 1 or 10 or 5 )$delta[1]  
}  
  
library("purr")  
1:5% > %map - db; (kf cv - err)
```

از درجه یک تا درجه ۵ را به ما می‌دهد.

Cross Validation for Classification Problem Miss Rate

فرق این روش در اینجا به این صورت هست که معیار ارزشیابی ما در اینجا دیگر MSE نخواهد بود و از نسبت بد طبقه بندی شده استفاده می کنیم.

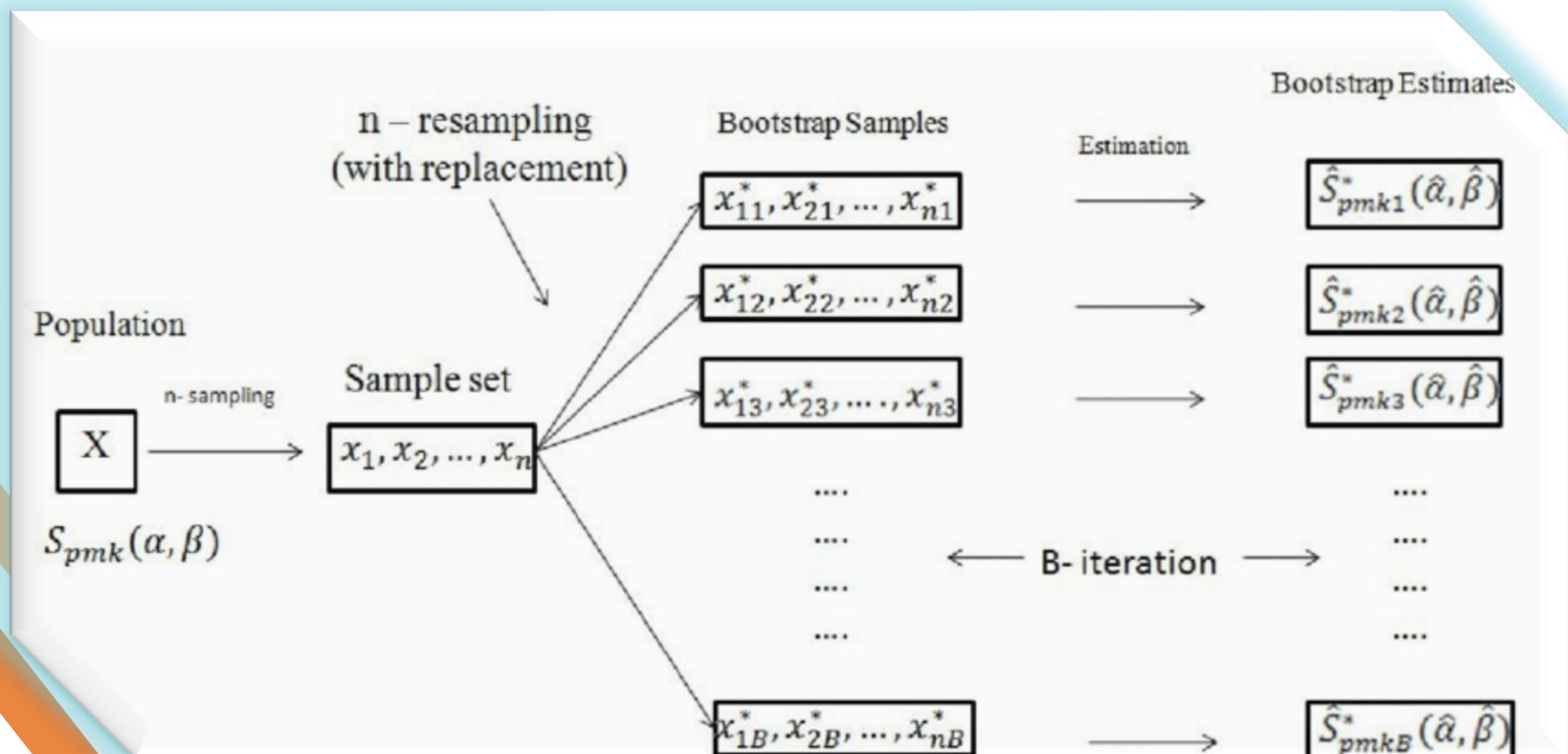
```
#R Codes, Cross Validation for Classification Problem Miss Rate  
glm.fit <- glm(y~x , family = "binomial", data = ... )  
Cost <- function(r, pi = 0){mean(abs(r - pi) > 0.5)}  
cv.glm(data, glm.fit, cost, k = 1 or k = 10 or k = 5)$delta[1]
```

مقدار خروجی بالا هرچه به ۵.۰ نزدیک تر باشد نشانه دقیق تر بودن پیشگویی مدل ما می باشد.

Boot strapping

خودگردان سازی:

یک روشی است که در آن عدم قطعیت یک برآوردگر را مورد بررسی قرار می‌دهیم. ما در اینجا از این روش برای بررسی ضرایب رگرسیونی خود استفاده می‌کنیم و زمانی که نتوانیم خطای استاندارد شده را برآورد بکنیم از این روش استفاده می‌کنیم.



Training and Test Errors

خطای آموزش:

خطای آموزش ما میزان خطایی هست که ما به کمک داده های آموزش خود می آییم و طبق معیار مورد نظر میزان دقت و خطای مدلی که بر اساس همین داده ها تشکیل شده است با اندازه گیری می کنیم.

خطای آزمون:

خطای آزمون ما میزان خطایی هست که ما به کمک داده های آزمون و یا تست خود روی مدلی که به کمک داده های آموزش تشکیل شده است تست کرده و میزان خطا را اندازه گیری می کنیم.

نکته: هدف ما کم بودن خطای آزمون هست و گاهی ممکن است خطای آموزش خیلی کم باشد ولی خطای آزمون بسیار زیاد.

Boot Stapping

#R Codes for Boot strapping Method

ایده کار در نرم افزار R برای این روش بدین صورت است که

(1) ساختن یک تابعی که آماره موردنظر (ضرایب رگرسیونی) را بسازد.

(2) تابع $boot$ را از پکیج $boot$ استخراج بکنیم تا بتوانیم از آن استفاده بکنیم.

```
library("boot")
Statistic <- function(data, index){
  glm.fit
  <- glm(y~poly(x, deg), data = ..., subset = index, family
  = ..., ...)
  coef(glm.fit)
}
statistic(data, index)
boot(data, statistic, 1000)
```

Criteria for Step3

برای انتخاب کردن معیارهای ارزیابی ۲ روش کلی وجود دارد.

1. بصورت غیرمستقیم خطای آزمون را محاسبه بکنیم و برآورد بکنیم.
2. بصورت مستقیم به بررسی خطای آزمون بپردازیم.

Evaluate Models

با توجه به C_p که در اسلاید قبلی درمورد آن بحث شد، می‌خواهیم معیارهای دیگر را معرفی بکنیم:

$$AIC = \frac{1}{\sigma^2} C_p$$

این معیار نیز، هرچه مقدار کمتری داشته باشد نشانه بهتر بودن مدل ما خواهد بود.

$$BIC = \frac{1}{n} (RSS + \log(n) d\hat{\sigma}^2)$$

این معیار نیز، هرچه مقدار کمتری داشته باشد نشانه بهتر بودن مدل ما خواهد بود، و اگر $\log(n) > 2$ یا $n > 7$ باشد، ضریب $d\hat{\sigma}^2$ نسبت به ضریب C_p بیشتر خواهد بود.

$$R^2 = 1 - \left(\frac{RSS}{SST} \right), \text{Adj. } R^2 = 1 - \left(\frac{MSE}{MST} \right)$$

این معیار نیز، هرچه مقدار بیشتری داشته باشد نشانه بهتر بودن مدل ما خواهد بود.

#R codes For AIC, BIC, C_p , Adjusted R squared:

result < -summary(best - subset)

Cp = result\$adjr2

BIC = result\$bic

which.max(adj - R2); which.min(Cp); which.min(BIC). اندیس را به ما می‌دهد.

روش مستقیم

در این روش ما بصورت مستقیم به بررسی خطای آزمون می‌پردازیم: در این روش ما باید خطای ارزیابی متقابل خود را برای مدل‌ها حساب کرده و بهترین را انتخاب کنیم. ما در اینجا برای مدل‌های مختلف با تعداد متغیرهای پیشگو مختلف یک خطای استاندارد (one standard error) داریم. و برای هر یک از این مدل‌های p متغیره، می‌آییم و مدلی را که کمترین میانگین انحراف استاندارد را دارد انتخاب می‌کنیم و آنرا با در نمودار خود $+$ و $-sd$ رسم می‌کنیم و بعد هر یک از مدل‌هایی که داخل این بازه قرار داشتند، ساده ترین آنها بهترین مدل ما خواهد بود.

$$X = \begin{matrix} obs1 \\ \vdots \\ obsn \end{matrix} \begin{bmatrix} 1 & x_1 & \dots & x_p \\ 1 & \vdots & \dots & \vdots \\ 1 & \vdots & \dots & \vdots \end{bmatrix}_{n \times p+1}$$

ما بعد از انتخاب بهترین مدل p متغیره خود، ضرایب آنرا استخراج می‌کنیم که آنها $p+1$ تا خواهند بود و بعد بصورت زیر مقادیر را پیش‌بینی خواهیم کرد.

$$\hat{y} = \beta X$$

روش مستقیم

```
#R codes for Cross Validation set Approches Method:  
Validation - erros <- vector("double",length=19)  
for(i in 1:19){coef - x <- coef(best - subset,id = i)  
pred - x <- test - m[,names(coef - x)]% * %coef - x  
validation - errors[i] = mean((test$y - pred - x)^2)}  
Plot(Validation - errors,type = "b")
```

روش مستقیم

#R codes for K Fold Cross Validation

```
predict.regsubset = function(object,newdata,id,...){  
  form <- as.formula(object$call[[2]])  
  mat <- model.matrix(form,newdata)  
  coefi <- coef(object,id = id)  
  xvars <- names(coefi)  
  mat[,xvars] % * % coefi  
  }  
  K <- 10  
  Folds <- sample(1:K , nrow(data), replace = TRUE)  
  Cv.errors <- matrix(NA,k,19,dimnames = list(NULL,paste(1:19)))  
  for(j in 1:k){  
    best_subset <- regsubset(y~x,data[folds == j,],nvmax = 19)  
    for(i in 1:19){  
      pred - x = predict.regsubset(best_subset,data[folds == j,],id = i)  
      Cv - errors[i,j] <- mean(data$y[folds == j] - pred - x^2)  
    }}
```

End

Thanks For Youre Attention

Join us ...