

# RISK ESTIMATION

## -QUANTIFYING THE WORLD-

Lecturer: Darren Homrighausen, PhD

# LOSS FUNCTIONS AND RISK

Define a function<sup>1</sup>  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  such that smaller values of  $\ell$  indicate **better** performance

Two important examples:

- $\ell(\hat{f}(X), Y) = (\hat{f}(X) - Y)^2$  (regression, square-error)
- $\ell(\hat{f}(X), Y) = \mathbf{1}(\hat{f}(X) \neq Y)$  (classification, 0-1)

These expressions are both **random variables**.

This leads us to define the prediction or estimation **risk** of a procedure  $\hat{f}$  to be

$$R(\hat{f}) = \mathbb{E}\ell(\hat{f}(X), Y)$$

---

<sup>1</sup>This is the loss for **prediction**. Other tasks, such as estimation, may have a different domain.

# RISK ESTIMATION

The prediction risk can be written

$$R(f) = \mathbb{E}\ell(f(X), Y) \leftrightarrow \text{Bias} + \text{Variance}$$

The overriding theme is that we would like to add a judicious amount of bias to get **lower** risk

As  $R$  isn't known, we need to estimate it

From the Preliminary Materials,  $\hat{R} = \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$  isn't very good

(In fact, one tends to not add bias when estimating  $R$  with  $\hat{R}$ )

# RISK ESTIMATION: A GENERAL FORM

The problem is that  $\hat{R}$  is overly **optimistic**

The **average optimism** is

$$\text{opt} =^* \mathbb{E}[R - \hat{R}]$$

Typically,  $\text{opt}$  is positive as  $\hat{R}$  will underestimate the risk

(\* See Elements of Statistical Learning, Chapter 7 for details for a more precise statement)

## RISK ESTIMATION: A GENERAL FORM

It turns out for a variety of  $\ell$  (such as squared error and 0-1)

$$\text{opt} = \frac{2}{n} \sum_{i=1}^n \text{Cov}(\hat{f}(X_i), Y_i)$$

This is related intimately with **degrees of freedom**

$$\text{df} = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(\hat{f}(X_i), Y_i) = \frac{n}{2\sigma^2} \text{opt}$$

$$(\sigma^2 = \mathbb{V}Y_i)$$

**EXAMPLE:** For multiple regression (i.e.  $\hat{f}(X) = \hat{\beta}_{LS}^\top X$ ),

$$\text{df} = \text{trace}(\mathbb{X}(\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top) = \text{rank}(\mathbb{X})$$

# A RISK ESTIMATE

Therefore, we get the following expression of risk

$$\text{GIC} = \hat{R} + \widehat{\text{opt}}$$

(Writing GIC indicates **generalized information criterion**)

Differing  $\widehat{\text{opt}}$  leads to AIC, BIC, Mallows Cp, and others

$\widehat{\text{opt}}$  depends on:

- a variance estimator  $\hat{\sigma}$
- a scaling term

# VARIOUS FORMS OF RISK ESTIMATES

$$\text{AIC} = \hat{R} + 2 \cdot \text{df} \cdot \hat{\sigma}^2/n$$

$$\text{AICc}(\hat{\beta}) = \text{AIC} + \frac{2\text{df} \cdot (\text{df} + 1)}{n - \text{df} - 1}$$

$$\text{BIC} = \hat{R} + \log(n) \cdot \text{df} \cdot \hat{\sigma}^2/n$$

Including more parameters leads to:

- a smaller  $\hat{R}$
- a larger  $\widehat{\text{opt}}$

**GOAL:** Now, we can use one of the GIC procedures to tell us which model to use

(As long as  $\log n \geq 2$ , BIC picks a **smaller** model than AIC)

# Cross-validation



# A DIFFERENT APPROACH TO RISK ESTIMATION

Let  $(X_0, Y_0)$  be a test observation, identically distributed as an element in  $\mathcal{D}$ , but also **independent** of  $\mathcal{D}$ .

$$R(f) = \mathbb{E} \ell(f(X_0), Y_0) \underbrace{=}_{\text{regression}} \mathbb{E} (Y_0 - f(X_0))^2$$

Of course, the quantity  $(Y_0 - f(X_0))^2$  is an unbiased estimator of  $R(f)$  and hence we could use it to estimate  $R(f)$

However, **we don't have any such new observation**

Or do we?

# AN INTUITIVE IDEA

Let's set aside one observation and predict it

**For example:** Set aside  $(X_1, Y_1)$  and fit  $\hat{f}^{(1)}$  on  $(X_2, Y_2), \dots, (X_n, Y_n)$

(The notation  $\hat{f}^{(1)}$  just symbolizes leaving out the first observation before fitting  $\hat{f}$ )

$$R_1(\hat{f}^{(1)}) = (Y_1 - \hat{f}^{(1)}(X_1))^2$$

As the left off data point is **independent** of the data points used for estimation,

$$\mathbb{E}R_1(\hat{f}^{(1)}) \approx R(\hat{f})$$

# LEAVE-ONE-OUT CROSS-VALIDATION

Cycling over all observations and taking the average produces  
leave-one-out cross-validation

$$\text{CV}_n(\hat{f}) = \frac{1}{n} \sum_{i=1}^n R_i(\hat{f}^{(i)}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}^{(i)}(X_i))^2.$$

# MORE GENERAL CROSS-VALIDATION SCHEMES

Let  $\mathcal{N} = \{1, \dots, n\}$  be the index set for  $\mathcal{D}$

- **K-FOLD:** Fix  $V = \{v_1, \dots, v_K\}$  such that  $v_j \cap v_k = \emptyset$  and  $\bigcup_j v_j = \mathcal{N}$

$$\text{CV}_K(\hat{f}) = \frac{1}{K} \sum_{v \in V} \frac{1}{|v|} \sum_{i \in v} (Y_i - \hat{f}^{(v)}(X_i))^2$$

- **BOOTSTRAP:** Instead of partitioning, we could make  $K$  bootstrap draws and average
- **FACTORIAL:** We could make all subsets of  $\mathcal{N}$  and average

# MORE GENERAL CROSS-VALIDATION SCHEMES: A COMPARISON

- $CV_K$  gets more computationally demanding as  $K \rightarrow n$
- The bias of  $CV_K$  goes down, but the variance increases as  $K \rightarrow n$
- The factorial version isn't commonly used  
(Very computationally demanding)

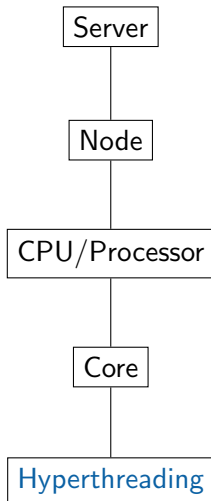
# SUMMARY TIME

- CV** Prediction risk consistent. Generally selects a model larger than necessary
- AIC** Provides optimal risk estimation (and is asymptotically equivalent to CV). Inconsistent for model selection
- BIC** Consistent for model selection consistent, sub optimal for risk estimation

Aside: There exist impossibility theorems stating that risk estimation procedures good at prediction are bad at model selection (and vice-versa)

# Parallelism

# DISTRIBUTED COMPUTING HIERARCHY



**EXAMPLE:** A server might have

- 64 nodes
- 2 processors per node
- 16 cores per processor
- **hyper threading**

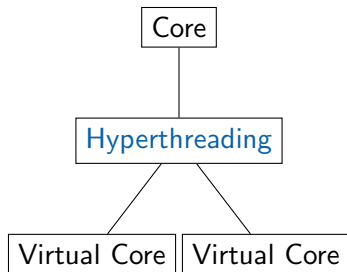
The goal is to somehow allocate a **job** so that these resources are used efficiently

Jobs are composed of **threads**, which are specific computations



# HYPERTHREADING

Developed by Intel, Hyperthreading allows for each core to pretend to be two cores



This works by trading off computation and read-time for each core

# GBM: FIGURES

```

train = agaricus.train
test = agaricus.test
nround = 10000
bst = xgboost(data = train$data, label = train$label, max.depth = 2,
              eta = 1, nthread = 1, nround = nround, objective =
"binary:logistic",
              print.every.n = nround/10)
pred = predict(bst, test$data)

##
# GBM
##
require(gbm)
data(agaricus.train, package='xgboost')
data(agaricus.test, package='xgboost')
Xtrain = data.frame(as.matrix(agaricus.train$data))
Ytrain = agaricus.train$label
test = agaricus.test
nround = 10000
gbml <-
gbm(Ytrain~.,data=Xtrain,      # formula
    distribution="bernoulli",  # see the help for other choices
    n.trees=200,              # number of trees
    shrinkage=0.1,            # shrinkage or learning rate,
                              # 0.001 to 0.1 usually work
                              # 1: additive model, 2: two-way interactions,
                              # etc.
    interaction.depth=3,
    # bag.fraction = 0.5,      # subsampling fraction, 0.5 is probably best
    # train.fraction = .9,     # fraction of data for training,
                              # first train.fraction*N used for training
    cv.folds = 3,              # minimum total weight needed in each node
    minobsinnode = 10,         # do 3-fold cross-validation
    keep.data=TRUE,           # keep a copy of the dataset with the object
    verbose=TRUE,              # don't print out progress
    n.cores=2)                 # use only a single core (detecting #cores is
                              # error-prone, so avoid here)

```

# 0.001 to 0.1 usually

# 1: additive model, 2:

# subsampling fraction,

bably best

in.fraction = .9,

# fraction of data for

# first train.fraction\*N

aining

nobsinnode = 10,

each node

olds = 3,

# do 3-fold cross

.data=TRUE,

th the object

ose=TRUE,

# don't print out

res=2)

# use only a sin

#cores is

inDeviance	ValidDeviance	StepSize	Improve
1.2169	nan	0.1000	0.0837
1.0823	nan	0.1000	0.0670
0.9669	nan	0.1000	0.0578
0.8699	nan	0.1000	0.0486
0.7772	nan	0.1000	0.0465
0.7062	nan	0.1000	0.0354
0.6434	nan	0.1000	0.0311

a good

formula = formula(data), distribution = "bernoulli", data = test, weights, var.monotone = NULL, n.trees = 100, interaction.depth = 1, minobsinnode = 10, shrinkage = 0.1, bag.fraction = 0.5, train.fraction = 1, cv.folds = 0, keep.data = TRUE, verbose = "cv", class.strategy = NA, n.cores = NA, 1)

Lecture 16 write up.pdf    scribe 14.pdf

On error: [c] ignore, [d] ignore all, [x] complete, [i] text to insert replacement

Show All

# PARALLELISM FOR CV IN R

Go to `crossValidationParallel.R`