

LINEAR METHODS FOR REGRESSION: RISK ESTIMATION

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

SUBSET SELECTION AND REGULARIZATION

For now, let's assume we are doing ordinary least squares, and hence the design (feature) matrix is $\mathbb{X} \in \mathbb{R}^{n \times p}$.

We want to do model selection for at least three reasons:

- **PREDICTION ACCURACY:** Can essentially *always* be improved by introducing some bias
- **INTERPRETATION:** A large number of features can sometimes be distilled into a smaller number that comprise the “big (little?) picture”
- **COMPUTATION:** A large p can create a huge computational bottleneck.

SUBSET SELECTION AND REGULARIZATION

We will address three related ideas

- **MODEL SELECTION:** Selection of only some of the original p features
- **DIMENSION REDUCTION/EXPANSION:** Creation of new features to help with prediction
- **REGULARIZATION:** Add constraints to optimization problems to provide stabilization

RISK ESTIMATION

REMINDER: Prediction risk is

$$R(f) = \mathbb{P}_{Z, \mathcal{D}} \ell_f \leftrightarrow \text{Bias} + \text{Variance}$$

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

As discussed, $\hat{R}_{\text{train}} = \hat{\mathbb{P}} \ell_f$ isn't very good

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

\hat{R}_{train} tends to **underestimate** R , hence we can call it **optimistic**

RISK ESTIMATION: A GENERAL FORM

Assume that we get a new draw of the training data, \mathcal{D}^0 , such that $\mathcal{D} \sim \mathcal{D}^0$ and

$$\mathcal{D} = \{(X_1, Y_1), \dots, (X_n, Y_n)\} \quad \text{and} \quad \mathcal{D}^0 = \{(X_1, Y_1^0), \dots, (X_n, Y_n^0)\}$$

If we make a small compromise to risk, we can form a sensible suite of risk estimators

To wit, letting $Y^0 = (Y_1^0, \dots, Y_n^0)^\top$, define

$$\begin{aligned} \underline{R}_{in} &= \mathbb{E}_{Y^0 | \mathcal{D}} \underbrace{\hat{\mathbb{P}}_{\mathcal{D}^0} \ell_{\hat{f}}}_{\ell_{\hat{f}} = \frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - Y_i^0)^2} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y^0 | \mathcal{D}} \ell(\hat{f}(X_i), Y_i^0) \\ &= \frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - Y_i^0)^2 = \frac{1}{n} \|\hat{f}(x) - Y^0\|_2^2 \end{aligned}$$

Then the **average optimism** is

$$\boxed{\text{opt}} = \mathbb{E}_Y [R_{in} - \hat{R}_{\text{train}}]$$

Typically, opt is positive as \hat{R}_{train} will underestimate the risk \rightarrow

RISK ESTIMATION: A GENERAL FORM

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

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

Therefore, we get the following expression of risk

$$\mathbb{E}_Y \underline{R}_{in} = \mathbb{E}_Y \hat{R}_{\text{train}} + \frac{2}{n} \sum_{i=1}^n \text{Cov}(\hat{f}(X_i), Y_i),$$

which has unbiased estimator (i.e. $\mathbb{E}_Y R_{\text{gic}} = \mathbb{E}_Y R_{in}$)

$$R_{\text{gic}} = \hat{R}_{\text{train}} + \frac{2}{n} \sum_{i=1}^n \text{Cov}(\hat{f}(X_i), Y_i)$$

DEGREES OF FREEDOM

We call the term (where $\sigma^2 = \mathbb{V} Y_i$)

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

the **degrees of freedom**

(This is really the **effective number of parameters**, with some caveats)

Our task now is to either estimate or compute opt to produce $\widehat{\text{opt}}$ and form:

$$\hat{R}_{\text{gic}} = \hat{R}_{\text{train}} + \widehat{\text{opt}}$$

This leads to Mallows Cp/Stein's unbiased risk estimator (SURE), as well as forms for AIC, BIC, and others

DEGREES OF FREEDOM: EXAMPLE

Sometimes the df is exactly computable.

(In other cases, it needs to be estimated)

Look at least squares regression onto \mathbb{X} , with $\mathbb{V}Y_i = \sigma^2$

$$\begin{aligned} df &= \frac{1}{\sigma^2} \sum_{(i)}^{\wedge} \text{cov}[\hat{f}(X_i), Y_i] = \dots = \frac{1}{\sigma^2} \text{trace} \left(\mathbb{X}(\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{V}Y \right) \\ &= \text{trace}(\mathbb{H}) = p \end{aligned}$$

INFORMATION CRITERIA

Of course, this isn't the usual way to introduce/conceptualize information criteria

For me, thinking of the **training error** as overly **optimistic** and correcting for that optimism is conceptually appealing

For others, forming a metric¹ on probability measures is more appealing

Let's go over this now for completeness

¹It will turn out to be a psuedo-metric; a small detail ▶

Comparing probability measures

KULLBACK-LEIBLER

Suppose we have data Y that comes from the probability density function f .

What happens if we use the probability density function g instead?

EXAMPLE: Suppose $Y \sim N(\mu, \sigma^2) = f$. We want to predict a new Y_* , but we model it as $Y_* \sim N(\mu_*, \sigma^2) = g$

How **far** away are we? We can either compare μ to μ_* or Y to Y^*
(This is the approach taken via the **optimism**)

Or, we can compute how **far** f is from g
(**far** indicates we need a notion of distance)

KULLBACK-LEIBLER

One central idea is **Kullback-Leibler** discrepancy²

$$\begin{aligned} KL(f, g) &= \int \log \left(\frac{f(y)}{g(y)} \right) f(y) dy \\ &\propto - \int \log(g(y)) f(y) dy \quad (\text{ignore term without } g) \\ &= -\mathbb{P}_f[\log(g(Y))] \end{aligned}$$

This gives us a sense of the **loss** incurred by using g instead of f .

²This has many features of a distance, but is not a true distance as $KL(f, g) \neq KL(g, f)$.

KULLBACK-LEIBLER DISCREPANCY

Usually, g will depend on some parameters, call them θ

EXAMPLE: In regression, we can specify $f = N(\overbrace{X^\top \beta_*}, \sigma_y^2)$ for a fixed (true)³ β , and let $g_\theta = N(X^\top \beta, \sigma^2)$ over all $\theta \in \mathbb{R}^p \times \mathbb{R}^+$

$$\Theta = (\beta, \sigma^2)$$

As $KL(f, g_\theta) = -\mathbb{P}_f[\log(g_\theta(Y))]$, we minimize this over θ .

Again, \mathbb{P}_f is unknown, so we minimize $-\log(g_\theta(Y))$ instead. This is the maximum likelihood value

$$\hat{\theta}_{ML} = \arg \max_{\theta} g_\theta(Y)$$

³We actually don't need to assume things about a true model nor have it be nested in the alternative models.

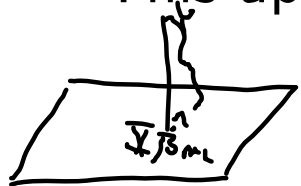
KULLBACK-LEIBLER DISCREPANCY

Now, to get an operational characterization of the KL divergence at the ML solution

$$-\mathbb{P}_f[\log(g_{\hat{\theta}_{ML}}(Y))]$$

we need an approximation (don't know f , still)

This approximation⁴ is exactly AIC:



$$\text{AIC} = -\log(g_{\hat{\theta}_{ML}}(Y)) + |\hat{\beta}_{ML}|$$

OF PARAMETERS
 $\beta_{ML} \rightarrow \hat{\beta}_{ML} = p$
 NORMAL σ
 $n \hat{R}_{train}$

Example:

Let $\log(g_{\theta}(y)) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|Y - \mathbb{X}\beta\|_2^2$

σ^2 KNOWN: $\hat{\beta} = \mathbb{X}^\dagger Y$

$$\text{AIC} \propto n\hat{R}_{train}/(2\sigma^2) + p = \hat{R}_{train} + 2\sigma^2 n^{-1} p$$

σ^2 UNKNOWN: $\hat{\beta} = \mathbb{X}^\dagger Y, n\hat{\sigma}^2 = (I - \mathbb{X}\mathbb{X}^\dagger)Y = n\hat{R}_{train}$

$$\text{AIC} \propto n \log(\hat{R}_{train})/2 + p = \log(\hat{R}_{train}) + 2n^{-1} p$$

$\frac{1}{2\hat{\sigma}_{ML}^2} \|Y - \mathbb{X}\hat{\beta}_{ML}\|_2^2$
 $\hat{\sigma}_{ML} = \frac{1}{n} \|Y - \mathbb{X}\hat{\beta}_{ML}\|_2^2$

⁴See "Multimodel Inference" Burnham, Anderson (2004)

SUMMARY

For \hat{R}_{gic} : OPTIMISM

$$\hat{R}_{\text{train}+\widehat{\text{opt}}} = \underbrace{\hat{R}_{\text{train}}}_{\text{Mallows Cp}} + 2\sigma^2 n^{-1} \text{df} = \begin{cases} \text{AIC, known } \sigma^2 \\ \text{Mallows Cp} \\ \text{SURE} \end{cases} \begin{array}{l} \text{if } \hat{f}(X) = X^\top \hat{\beta}_{LS} \\ \text{most } \hat{f}(X) \end{array}$$

Or

KL

$$\text{IC} = \log(\hat{R}_{\text{train}}) + c_n n^{-1} \text{df} = \begin{cases} \text{AIC, unknown } \sigma^2 & \text{if } c_n = 2 \\ \text{BIC} & \text{if } c_n = \log(n) \end{cases}$$

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} .

Prediction risk: $R(f) = \underbrace{\mathbb{E}}_{\text{UNKNOWN}} (Y_0 - f(X_0))^2$

Of course, the quantity $(Y_0 - f(X_0))^2$ is an unbiased estimator of $R(f)$ and hence we could estimate $R(f)$ $\mathbb{E}(Y_0 - f(X_0))^2 = R(f)$

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,

$$\mathbb{E}_{(X_1, Y_1) | \mathcal{D}_{(1)}} R_1(\hat{f}^{(1)}) \stackrel{D}{=} R(\hat{f}(\mathcal{D}_{n-1})) \approx R(\hat{f}(\mathcal{D}))$$

LEAVE-ONE-OUT CROSS-VALIDATION

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

LOOCV

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

MORE GENERAL CROSS-VALIDATION SCHEMES

Let $\mathcal{N} = \{1, \dots, n\}$ be the index set for $\mathcal{D} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$

Define a distribution \mathcal{V} over \mathcal{N} with (random) variable v

Then, we can form a general **cross-validation** estimator as

$$\mathcal{V}(i) = \frac{1}{n} \Rightarrow \text{LOO} \quad \text{CV}_{\mathcal{V}}(\hat{f}) = \mathbb{E}_{\mathcal{V}} \hat{\mathbb{P}}_{\mathcal{V}} \ell_{\hat{f}(v)}$$

$$\hat{\mathbb{P}}_{\mathcal{V}} f = \frac{1}{n} \sum_{i=1}^n f(X_i)$$

$$\mathcal{V} = \{i\}$$

$$\ell_f(x) = (Y - f(x))^2 \Rightarrow \hat{\mathbb{P}}_{\mathcal{V}} \ell_{\hat{f}(v)} = (Y_i - \hat{f}^{(i)}(X_i))^2$$

MORE GENERAL CROSS-VALIDATION SCHEMES: EXAMPLES

LOOKS \Rightarrow n -FOLD CV

$$CV_{\mathcal{V}}(\hat{f}) = \mathbb{E}_{\mathcal{V}} \hat{\mathbb{P}}_{\mathcal{V}} \ell_{\hat{f}^{(\mathcal{V})}}$$

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

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

- **BOOTSTRAP:** Let \mathcal{V} be given by the bootstrap distribution over \mathcal{N} (that is, sampling with replacement many times)
- **FACTORIAL:** Let \mathcal{V} be given by all subsets (or a subset of all subsets) of \mathcal{N} (that is, putting mass $1/(2^n - 2)$ on each subset)

MORE GENERAL CROSS-VALIDATION SCHEMES: A COMPARISON

$$\mathbb{E} CV_K = \frac{1}{K} \sum_{V \in \mathcal{V}} \frac{1}{|V|} \sum_{i \in V} \mathbb{E} (Y_i - \underbrace{\hat{f}^{(V)}(X_i)}_{n-|V|})^2$$

$$\mathbb{E} \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}^{(i)}(X_i))^2 = \frac{1}{n} \sum_{i=1}^n \mathbb{E} (Y_i - \underbrace{\hat{f}^{(i)}(X_i)}_{n-1})^2$$

- CV_K gets more computationally demanding as $K \rightarrow n$
- The bias of CV_K goes down, but the variance increases as $K \rightarrow n$ $K=10$
- The factorial version isn't commonly used except when doing a 'real' data example for a methods paper
- ~~There are many other flavors of CV. One of them, called "consistent cross validation" [HOMEWORK] is a recent addition that is designed to work with sparsifying algorithms~~

$$\mathbb{V} LOOCV = \frac{1}{n^2} \left(\sum W + \sum \sum CV \right) \rightarrow \text{VARIANCE IS LARGE}$$

Summary time

RISK ESTIMATION METHODS

- CV** Prediction risk consistent (Dudoit, van der Laan (2005)). Generally selects a model larger than necessary (unproven)
- AIC** Minimax optimal risk estimator (Yang, Barron (1998)). Model selection inconsistent*
- BIC** Model selection consistent (Shao (1997) [low dimensional]. Wang, Li, Leng (2009) [high dimensional]). Slow rate for risk estimation*

(Stone (1977) shows that CV_n and AIC are asymptotically equivalent.)

(*Yang (2005) gives an impossibility theorem: for a linear regression problem it is impossible for a model selection criterion to be both consistent and achieve minimax optimal risk estimation)