

STA 314: Statistical Methods for Machine Learning I

Lecture 3 - Model selection under linear model Cross-validation

Xin Bing

Department of Statistical Sciences
University of Toronto

- We have learned the OLS approach:

$$\hat{\beta} = \underset{\alpha}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\alpha\|_2^2.$$

- We have learned the statistical properties of $\hat{\beta}$ under the linear model

$$Y = \beta_0 + X_1\beta_1 + \cdots + X_p\beta_p + \epsilon.$$

- ▶ Unbiasedness
- ▶ Estimation error (ℓ_2 -norm)
- ▶ Inference (confidence intervals, hypothesis testing).

Why consider alternatives to the OLS estimator?

Alternative fitting procedures to OLS:

- **prediction / estimation**: the OLS estimator

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

has large variance when p is large. Especially, if $p > n$, then OLS estimator is not unique and its variance is infinite.

- **interpretability**: By removing irrelevant features – that is, by setting some coefficient estimates to zero – we can obtain a model that is more parsimonious hence more interpretable.

What are the alternatives?

1. **Subset Selection.** We identify a subset of the p predictors that we believe to be related to the response. We then fit a model using the OLS approach on the identified set of predictors.
2. **Shrinkage.** We fit a model involving all p predictors, but the estimated coefficients are shrunk towards zero relative to the OLS estimator. This shrinkage (also known as regularization) has the effect of reducing variance. Some could also perform variable selection.
3. **Dimension Reduction.** We project the p predictors into a M -dimensional subspace, where $M < p$. This is achieved by computing M different linear combinations, or projections, of the original predictors. Then the resulting M projections are used as new predictors to fit a linear regression model by OLS.

How to choose the optimal one among a set of models?

Example

$$\text{Model 1: } Y = \alpha_0 + \alpha_1 X_1 + \epsilon$$

$$\text{Model 2: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

lead to different predictors at $X = \mathbf{x} = (x_1, x_2)$

$$\hat{f}_1(x) = \hat{\alpha}_0 + \hat{\alpha}_1 x_1 \quad \text{v.s.} \quad \hat{f}_2(x) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2.$$

Ideally, we choose the one that has a **smaller expected MSE**.

How to compare expected MSEs?

- When we have \mathcal{D}_{test} , we compare the test MSE errors directly

$$\frac{1}{m} \sum_{i=1}^m \left(y_i^{(T)} - \hat{\alpha}_0 - \hat{\alpha}_1 x_{i1}^{(T)} \right)^2$$

v.s.

$$\frac{1}{m} \sum_{i=1}^m \left(y_i^{(T)} - \hat{\beta}_0 - \hat{\beta}_1 x_{i1}^{(T)} - \hat{\beta}_2 x_{i2}^{(T)} \right)^2$$

- What if we don't have \mathcal{D}_{test} ?

There are two common approaches for model selection when we don't have \mathcal{D}_{test} :

- We can directly estimate the expected MSE by manually creating a “test set” using data-splitting techniques:
 - ▶ validation set approach
 - ▶ cross-validation approach
- We can avoid estimating the expected MSE by making an adjustment to the training error to account for the model complexity:
 - ▶ Mallow's C_p
 - ▶ adjusted R^2
 - ▶ AIC & BIC

Direct estimation of the expected MSE: data-splitting techniques

We randomly split the available data to create a validation set that functions as a test set.

- Validation set approach: one-time data splitting
- Cross-validation approach: multiple-time data splitting

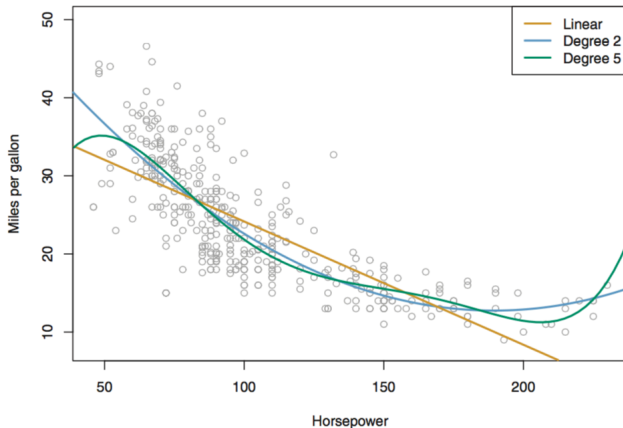
Validation set approach

- *Randomly* divide the available set of samples into two parts: a **training set** and a **validation** (hold-out) set.
 - ▶ What is the proportion? Depends.
- The model is fitted on the training set, and the fitted model is evaluated by the validation set.
- The resulting validation-set MSE provides an estimate of the expected MSE.



Example: Auto Data

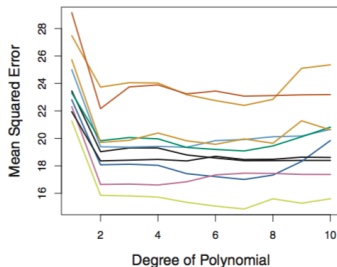
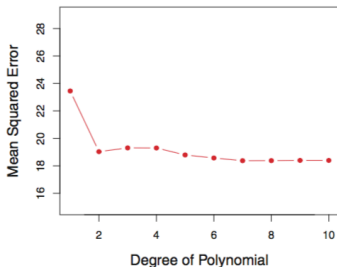
In Lecture 2, we find there appears to be a non-linear relationship between **mpg** and **horsepower**.



Whether a cubic or higher-order predictor provides a better fit?

Example: Auto Data – Compare linear vs higher-order polynomial terms in a linear regression.

We randomly split the 392 observations into two sets, a training set containing 196 of the data points, and a validation set containing the remaining 196 data.



- Left: Validation error estimates for a single split into training and validation data sets.
- Right: Validation method repeated 10 times with each time using a different random split of the observations into a training set and a validation set.
- We can see the one-time data splitting is not stable

Drawbacks of Validation Set Approach

- The validation estimate of the test error can be highly unstable, depending on which observations are included in the training set and which are in the validation set.
- Only a subset of the observations – those in the training set rather than in the validation set – are used to fit the model.
- This suggests that the validation set error may tend to overestimate the test error for the model fit on the entire data set.
- How to remedy these drawbacks?

Leave-One-Out Cross-Validation (LOOCV)

- First split the data into two parts by leaving out the **first** observation:
 - ▶ a validation set: (\mathbf{x}_1, y_1)
 - ▶ a training set: the remaining observations $(\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
 - ▶ using the training set, we fit the model \hat{f}_1 and predict y_1 as $\hat{f}_1(\mathbf{x}_1)$ using the value \mathbf{x}_1 . The test error could be approximated by

$$MSE_1 = (y_1 - \hat{f}_1(\mathbf{x}_1))^2.$$

- ▶ not good enough!

Leave-One-Out Cross-Validation

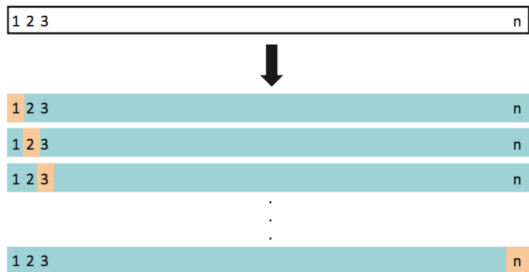
- Repeat the procedure by leaving out the **second** observation:
 - ▶ a validation set: (\mathbf{x}_2, y_2) ,
 - ▶ a training set: the remaining observations $(\mathbf{x}_1, y_1), (\mathbf{x}_3, y_3), \dots, (\mathbf{x}_n, y_n)$
 - ▶ using the training set, we fit the model \hat{f}_2 and predict y_2 as $\hat{f}_2(\mathbf{x}_2)$ using the value \mathbf{x}_2 . Compute

$$MSE_2 = (y_2 - \hat{f}_2(\mathbf{x}_2))^2.$$

- Repeating the approach n times by leaving out **each** observation to obtain MSE_1, \dots, MSE_n .
- The LOOCV estimate for the test MSE is the average of these n test error estimates:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i.$$

Leave-One-Out Cross-Validation



Validation data sets in beige, and training sets in cyan.

LOOCV vs Validation Set Approach

LOOCV has the following advantage over the validation set approach.

- The training set of LOOCV is almost the same as the entire data set. The fitted model is almost as good as that based on the entire data set.
- The validation approach yields different results when applied repeatedly, because the training/validation set is randomly divided. LOOCV has no randomness in the splitting.

However, LOOCV can be computationally expensive in general.¹

¹In linear model, the computation can be simplified, the formula is shown in page 202 of the textbook.

k-Fold Cross-Validation

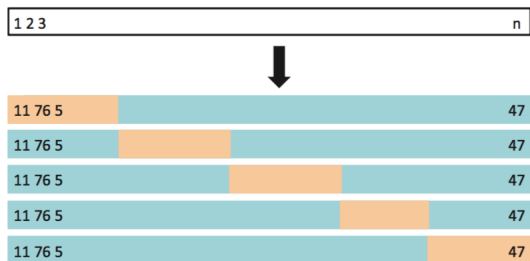
- **k-fold CV** is to randomly divide the data into k (roughly) equal-sized groups or folds.
- The first fold is treated as a validation set, and the method is fit on the remaining $k - 1$ folds. We compute the mean squared error, MSE_1 , for the observations in the first fold.
- Then we repeat the procedure to fold 2, fold 3, ..., fold k , and get $MSE_2, MSE_3, \dots, MSE_k$.
- The k-fold CV estimate is computed by averaging these values,

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i.$$

Remark

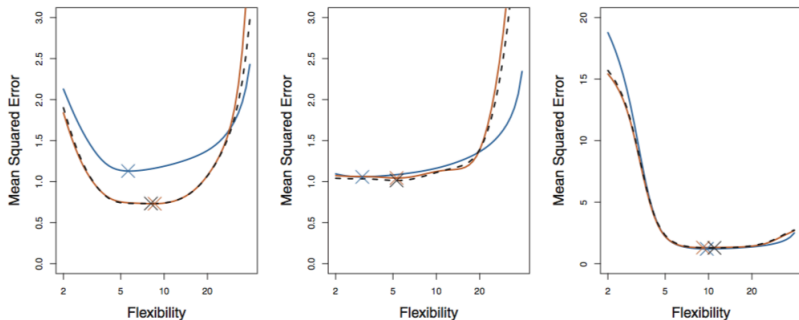
- LOOCV is a special case of n -fold CV.
- 5-fold or 10-fold is commonly used in practice.

k-Fold Cross-Validation



Validation data sets in beige, and training sets in blue.

k-Fold Cross-Validation



True test MSE (in blue), the LOOCV estimate (black dashed line), and the 10-fold CV estimate (in orange) for three simulated data sets.

Cross-Validation on Classification Problems

- Cross-validation also works for classification problems.
- For LOOCV, we split the data in the same way as before. We compute the error on the validation set

$$Err_1 = 1 \{y_1 \neq \hat{f}(x_1)\}.$$

- Then we repeat the procedure n times, and get $Err_2, Err_3, \dots, Err_n$.
- The LOOCV estimate is computed by averaging these values,

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n Err_n.$$

Cross-Validation is sometimes tricky!

Independence between the fitted model and the validation set is the key! Meaning that you should **NOT** use the validation set to fit your model.

Example (A tricky one)

Consider a simple two-step approach applied to some data \mathcal{D}^{train} .

- Starting with 5000 predictors and 100 samples, find the 10 predictors having the largest correlation with the outcome.
- We then apply the OLS using only these 10 predictors.

How do we estimate the expected MSE of the fitted model from this approach?

Avoid estimating the expected MSE:

C_p , AIC, BIC, and adjusted R^2

- These techniques adjust the training error for the model complexity.
- They are limited to
 - ▶ parametric models such as linear model
 - ▶ models where the data likelihood is correctly specified

Cont'd: Avoid estimating the expected MSE

For any given fitted model \hat{f} , let $\hat{f}(\mathbf{x}_i)$ be the fitted value for the i th observation. For instance, for a fitted linear model with p predictors,

$$\hat{f}(\mathbf{x}_i) = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_p x_{ip}.$$

Recall that $y_i - \hat{f}(\mathbf{x}_i)$ is the i th residual. The residual sum of squares (RSS) is defined as

$$\frac{1}{n} \text{RSS}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2.$$

Yes, this is indeed the training MSE of \hat{f} ! It gets smaller as p increases.

Let p be the total # of parameters in the model

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

When σ^2 is unknown, one use a consistent estimator $\hat{\sigma}^2$.

- C_p adds a penalty $2p\hat{\sigma}^2/n$ to the training MSE to adjust for the fact that the training error is always in favor of more complex models.
- we choose the model with the lowest C_p value.
- C_p is mainly for selecting linear predictors in regression.

Adjusted R^2

Recall that the total sum of squares (TSS) is defined as

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Recall that

$$R^2(\hat{f}) = \frac{TSS - RSS(\hat{f})}{TSS} = 1 - \frac{RSS(\hat{f})}{TSS}.$$

By contrast,

$$\text{Adjusted } R^2(\hat{f}) = 1 - \frac{RSS(\hat{f})/(n-p-1)}{TSS/(n-1)}.$$

Remark. Unlike the R^2 statistic, the adjusted R^2 statistic pays a price for the inclusion of unnecessary variables in the model.

Adjusted R^2 vs R^2

$$\text{Adjusted } R^2(\hat{f}) = 1 - \frac{RSS(\hat{f})/(n - p - 1)}{TSS/(n - 1)}.$$

- Note

$$\operatorname{argmax}_{\hat{f}} \text{Adjusted } R^2(\hat{f}) = \operatorname{argmin}_{\hat{f}} \frac{RSS(\hat{f})}{n - p - 1}.$$

While RSS always decreases as p increases, $RSS/(n - p - 1)$ may increase or decrease.

- A **larger** value of adjusted R^2 indicates a model with smaller test error.
- Both C_p and adjusted R^2 are restricted to selection of linear models.

Let \hat{f} be the fitted model obtained from the MLE approach so that $L(\hat{f})$ is the maximum of the likelihood function.

- **AIC:**

$$AIC(\hat{f}) = -2 \log L(\hat{f}) + 2p,$$

In the linear model with $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$, $AIC(\hat{f})$ is proportional to $C_p(\hat{f})$, selecting the same model.

- **BIC:**

$$BIC(\hat{f}) = -2 \log L(\hat{f}) + (\log n)p,$$

BIC places a heavier penalty $(\log n)p$ on models with many predictors, and hence results selecting smaller-size models than AIC and C_p .

- For both AIC and BIC, we select the best model that has the lowest value.
- To compute AIC and BIC, we need to specify the likelihood, i.e. the distribution of $Y | X$, and to compute the maximum likelihood estimator.
- AIC and BIC can also be used for selecting parametric models in classification problems.

Example (Revisited)

$$\text{Model 1: } Y = \alpha_0 + \alpha_1 X_1 + \epsilon$$

$$\text{Model 2: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

lead to different predictors at $X = x = (x_1, x_2)$

$$\hat{f}_1(x) = \hat{\alpha}_0 + \hat{\alpha}_1 x_1 \quad \text{v.s.} \quad \hat{f}_2(x) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2.$$

Their residual sum of squares (RSS) can be computed as:

$$RSS(\hat{f}_1) = \sum_{i=1}^n (y_i - \hat{\alpha}_0 - \hat{\alpha}_1 x_{i1})^2$$

$$RSS(\hat{f}_2) = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2})^2$$

- Mallow's C_p

$$C_p(\hat{f}_1) = \frac{1}{n} \left(\text{RSS}(\hat{f}_1) + (2 \times 1)\sigma^2 \right)$$

$$C_p(\hat{f}_2) = \frac{1}{n} \left(\text{RSS}(\hat{f}_2) + (2 \times 2)\sigma^2 \right).$$

- Adjusted R^2

$$\text{Adjusted } R^2(\hat{f}_1) = 1 - \frac{\text{RSS}(\hat{f}_1)/(n-2)}{\text{TSS}/(n-1)}$$

$$\text{Adjusted } R^2(\hat{f}_2) = 1 - \frac{\text{RSS}(\hat{f}_2)/(n-3)}{\text{TSS}/(n-1)}.$$

Discussion / recommendation on these two approaches

- The data-splitting technique has two advantages relative to C_p , adjusted R^2 , AIC and BIC, :
 - ▶ it provides a direct estimate of the test error
 - ▶ It can also be used in a wider range of model selection tasks, even in cases where it is hard to pinpoint the model degrees of freedom (e.g. the number of predictors in the model) or hard to estimate the error variance.
- The data-splitting technique also has a couple of drawbacks comparing to the other approach:
 - ▶ it requires a relatively large sample size
 - ▶ it is difficult to have guarantees for the model selected by using CV.
 - ▶ when the distribution is specified and the error of variance can be consistently estimated, the first approach is preferred.

What's next?

Apply these two techniques
for model selection under linear models.

Alternatives to the OLS using all predictors:

- **Subset Selection.** Identify a subset of the p predictors that we believe to be related to the response. Then fit the model by using the identified predictors via OLS.
 - ▶ Best Subset Selection
 - ▶ Stepwise Selection
- **Shrinkage Regression**
 - ▶ Ridge
 - ▶ Lasso
- **Dimension Reduction.** Later after PCA.

Example

Suppose we have access to i.i.d. samples of the response Y and the features

$$X = (X_1, X_2, X_3).$$

For fitting a regression that is linear in X , what are the all possible subsets?

Algorithm 6.1 *Best subset selection*

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
 2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

- Step 2 identifies the best model for each subset size. RSS can be used here. Why?
- In Step 3, can we use *RSS* or R^2 ?

Best Subset Selection

- The same approach can be used for other types of models, such as logistic regression (RSS replaced by deviance).
- However! For best subset selection, we need to fit and compare

$$\binom{p}{0} + \binom{p}{1} + \binom{p}{2} + \cdots + \binom{p}{p} = 2^p$$

models.

Example (Revisited)

Suppose we have access to i.i.d. samples of the response Y and the features

$$X = (X_1, X_2, X_3).$$

What are the models we consider for forward stepwise?

Algorithm 6.2 *Forward stepwise selection*

1. Let \mathcal{M}_0 denote the *null* model, which contains no predictors.
 2. For $k = 0, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Forward Stepwise Selection

- **Pros:** It has computational advantage over best subset selection. In the k th iteration, we fit and compare $(p - k)$ models. In total,

$$1 + \sum_{k=0}^{p-1} (p - k) = 1 + \frac{p(p + 1)}{2}$$

models are considered, *much fewer* than 2^p models.

- **Cons:** It is a greedy procedure!
So not guaranteed to find the best possible model out of all 2^p models containing subsets of the p predictors.

The Credit Card Data

# Variables	Best subset	Forward stepwise
One	<code>rating</code>	<code>rating</code>
Two	<code>rating, income</code>	<code>rating, income</code>
Three	<code>rating, income, student</code>	<code>rating, income, student</code>
Four	<code>cards, income, student, limit</code>	<code>rating, income, student, limit</code>

Example (Revisited)

Suppose we have access to i.i.d. samples of the response Y and the features

$$X = (X_1, X_2, X_3).$$

What are the models we consider for backward stepwise?

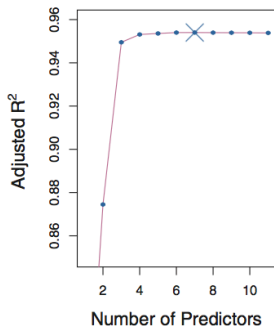
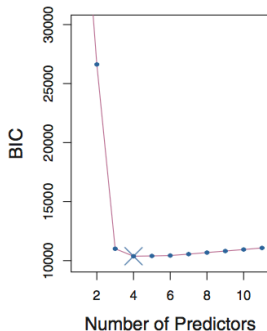
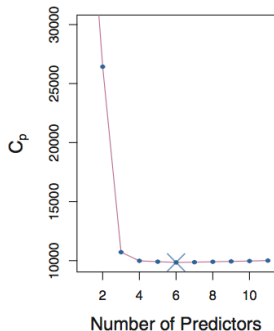
Algorithm 6.3 *Backward stepwise selection*

1. Let \mathcal{M}_p denote the *full* model, which contains all p predictors.
 2. For $k = p, p - 1, \dots, 1$:
 - (a) Consider all k models that contain all but one of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

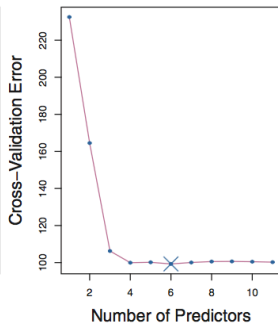
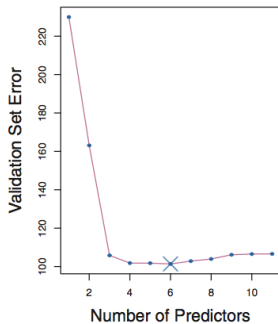
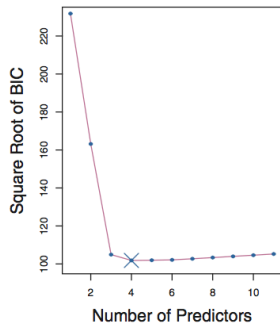
Backward Stepwise Selection

- For backward stepwise selection, we also compare $1 + p(p + 1)/2$ models, much fewer than 2^p models.
- Still a greedy approach!
It is not guaranteed to find the best possible model out of all 2^p models containing subsets of the p predictors.

The Credit Card Data: best subset selection via Mallows's C_p , BIC and adjusted R^2



The Credit Card Data: model selection via sample-splitting



Summary on subset selection

- Best subset selection will select the best model, as long as computation is affordable.
- Forward / Backward stepwise selection is computationally fast, but is not guaranteed to find the best model.
- What should we do in practice?