

Linear model selection and regularization

SDS 323

James Scott (UT-Austin)

Reference: Introduction to Statistical Learning Chapter 6

Outline

1. Goals of model selection and regularization
2. Model selection by stepwise selection
3. Regularized regression: an overview
4. Cross validation
5. Aside on software and sparse matrices
6. Dimension reduction

Note: although we talk about *regression* here, everything applies to logistic regression as well (and hence classification).

Model selection and regularization

In this set of notes, we're still talking about the good ol' linear model:

$$E(y \mid x) = \beta_0 + \sum_{j=1}^p \beta_j x_{ij}$$

But we're focused on improving the linear model by using some other model-fitting strategy beyond ordinary least squares (OLS).

Why move beyond OLS?

1. Improved prediction accuracy—sometimes radically improved.
2. More interpretable models.

Reason 1: improved prediction accuracy

- A linear model is generally thought of as a “high bias, low variance” kind of estimator, at least compared with alternatives we've seen (and others we have yet to see).
- But if the number of observations n is not much larger than the number of features p , even OLS can have high estimation variance.
- The result: overfitting and poor prediction accuracy.
- One solution: *shrinking* or *regularizing* the coefficient estimates so that they don't just chase noise in the training data.
- Extreme case: if $p > n$, there isn't even a unique OLS solution! No option but to do something else.

Reason 2: model interpretability

- A major source of variance in OLS estimation is the inclusion of unnecessary variables in the regression equation. (Here “unnecessary” means “ $\beta_j \approx 0$ ”.)
- This situation is especially pronounced in situations of *sparsity*: where you have lots of candidate features for predicting y , only a few of them are actually useful, but you don't know which ones.
- By removing these variables (i.e. setting their coefficients to zero), we get a simpler model that is easier to interpret and communicate.
- In this chapter, we'll learn some ways for doing this automatically.

Model selection

The seemingly obvious approach is *exhaustive enumeration*:

- fit all possible models under to a training set
- measure the generalization error of each one on a testing set.

But this can be *too* exhausting. What are the limiting performance factors here?

- it might take a long time to fit each model (when?)
- if so, it will take even longer to repeatedly re-fit each model to multiple training sets
- and there might be way too many models to check them all.

Iterative model selection

Thus a more practical approach to model-building is *iterative*:

- Start with a *working model*.
- Consider a limited set of “small” changes to the working model.
- Measure the performance gains/losses resulting from each small change.
- Choose the “best” small change to the working model. This becomes the new working model.
- Repeat until you can't make any more changes that make the model better.

Iterative model selection

OK, so:

- Where do I start?
- What is a “small change” to a model?
- What is a “limited set” of such changes?
- How do we measure performance gains/losses?

Forward selection

Suppose we have p candidate variables and interactions (called the “scope”). Start with a working model having no variables in it (the null model).

1. Consider all possible one-variable additions to the working model, and measure how much each addition improves the model's performance. (Note: we'll define “improvement” here soon!)
2. Choose the single addition that improves the model the most. This becomes the new working model. (If no addition yields an improvement, stop the algorithm.)
3. Repeat steps 1 and 2.

Backward selection

Suppose we have p candidate variables and interactions. Start with a working model having all these variables in it (the *full* model).

1. Consider all possible one-variable deletions to the working model, and measure how much each deletion improves the model's performance.
2. Choose the single deletion that improves the model the most. This becomes the new working model. (If no deletion yields an improvement, stop the algorithm.)
3. Repeat steps 1 and 2.

Stepwise selection (forward/backward)

Suppose we have p candidate variables and interactions (the scope). Start with any working model containing some subset of these variables. Ideally this should be a reasonable guess at a good model.

1. Consider all possible one-variable additions or deletions to the working model, and measure how much each addition or deletion improves the model's performance.
2. Choose the single addition or deletion that improves the model the most. This becomes the new working model. (If no deletion yields an improvement, stop the algorithm.)
3. Repeat steps 1 and 2.

Measuring model performance

- Ideally, we'd measure model performance using out-of-sample MSE, calculated by cross-validation.
- But this adds another computationally expensive step to the algorithm.
- Thus stepwise selection usually involves some approximation.
- Thus most statistical software will measure performance *not* by actually calculating MSE_{out} on some test data, but rather using one of several possible heuristic approximations for MSE_{out} .

AIC

The most common one is called AIC (“Akaike information criterion”):

$$\text{MSE}_{\text{AIC}} = \text{MSE}_{\text{in}} \left(1 + \frac{p}{n} \right) ,$$

where

$$\text{MSE}_{\text{in}} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - p}$$

is the usual unbiased estimated of the residual variance σ^2 in a linear model. Thus the AIC estimate is an *inflated* version of the in-sample MSE. **It is not a true out-of-sample estimate.**

AIC: a technical note

The general definition of AIC for a model with p coefficients is

$$\text{AIC} = \text{Deviance} + 2 \cdot p$$

where you'll recall that the deviance is -2 times the log likelihood of the model.

This is an estimate of the out-of-sample deviance of a model. Note that we inflate the in-sample (fitted) deviance by an additive factor of $2 \cdot P$.

In the special case of a linear model fit by OLS, this general definition reduces to the version of MSE_{AIC} just given.

AIC: some notes

- The inflation factor of $(1 + p/n)$ is always larger than 1.
- But the more parameters p you have relative to data points n , the larger the inflation factor gets.
- The AIC estimate of MSE is just an approximation to a true out-of-sample estimate. It still relies upon some pretty specific mathematical assumptions (linear model true, error Gaussian) that can easily be wrong in practice.

AIC: some notes

- Don't view AIC + stepwise selection algorithm as having God-like powers for discerning the single best model
- Don't treat it as an excuse to be careless. Instead proceed cautiously. Always verify that the stepwise-selected model makes sense and doesn't violate any crucial assumptions.
- It's also a good idea to perform a quick train/test split of your data and compute MSE_{out} for your final model.
- This ensures you're actually improving the generalization error versus your baseline model.

Problems with subset/stepwise selection

- `step()` is very slow.
- This is a generic problem with stepwise selection: each candidate model along the path must be refit from scratch.
- A related subtle (but massively important) issue is stability. MLEs can have high sampling variability where $p \approx n$: they change a lot from one dataset to another. So which MLE model is “best”“ changes a lot.
- AIC is only an *estimate* of out-of-sample performance. Sometimes a bad one! Stepwise is *too expensive* to do true out-of-sample validation.

Regularization

The key to modern statistical learning is **regularization**: departing from optimality to stabilize a system.

This is very common in engineering:

- Would you drive on an “optimal” bridge (absolute minimum concrete required to support a load)?
- Would you fly on an “optimal” airplane? (Ask Boeing!)

Regularization: some intuition

Recall that deviance ($= -2 \cdot LHD$) is the cost of mis-fitting the data:

- It penalizes discrepancy between model predictions \hat{y} and actual data y .
- This cost is *explicit* in maximum likelihood estimation.

But nonzero choices of β_j also have costs:

- $\beta_j = 0$ is a “safe” (zero-variance) choice.
- Any other nonzero choice of β_j incurs a cost: the need to use data (a finite resource) to estimate it. You pay this cost in **variance**.
- This cost is *hidden* in maximum likelihood estimation.
- Solution: make the cost explicit!

Regularization

The “optimal” fit $\hat{\beta}_{\text{MLE}}$ maximizes the likelihood, or equivalently minimizes the deviance, as a function of β :

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \text{dev}(\beta)$$

The regularized fit minimizes the deviance *plus a “penalty”* on the complexity of the estimate:

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \text{dev}(\beta) + \lambda \cdot \text{pen}(\beta)$$

Here λ is the penalty weight, while “pen” is some cost function that penalizes departures of the fitted β from 0.

Regularization: two common penalties

Ridge regression:

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \text{dev}(\beta) + \lambda \cdot \sum_{j=1}^p \beta_j^2$$

Lasso regression:

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \text{dev}(\beta) + \lambda \cdot \sum_{j=1}^p |\beta_j|$$

The advantage of the lasso approach is that it gives sparse solutions: many of the $\hat{\beta}_j$'s are identically zero. This yields *automatic variable selection*.

Regularization: path estimation

The lasso fits $\hat{\beta}$ to minimize $\text{dev}(\beta) + \lambda \cdot \sum_{j=1}^p |\beta_j|$. We'll do this for a sequence of penalties $\lambda_1 > \lambda_2 > \dots > \lambda_T$.

We can then apply standard model-selection tools (e.g. out-of-sample validation, AIC) to pick the best λ .

Path estimation:

- Start with big λ_1 : so big that $\hat{\beta} = 0$.
- For $t = 2, \dots, T$, update $\hat{\beta}$ to be optimal under $\lambda_t < \lambda_{t-1}$.

Regularization: path estimation

Some key facts:

- The estimate $\hat{\beta}$ changes continuously along this path of λ 's.
- It's fast! Each update is computationally very easy.
- It's stable: the optimal λ may change a bit from sample to sample, but these changes don't affect the overall fit that much.

This is like a better version of traditional stepwise selection. **But how do we choose an optimal λ ?**

Recall train/test splits

Goal: estimate the prediction error of a fitted model.

Solution so far:

- Randomly split the data into a training set and testing set.
- Fit the model(s) on the training set.
- Make predictions on the test set using the fitted model(s).
- Check how well you did (RMSE, classification error, deviance...)

A perfectly sensible approach! But still has some drawbacks...

Train/test splits: drawbacks

Drawback 1: the estimate of the error rate can be highly variable.

- It depends strongly on which observations end up in the training and testing sets.
- Can be mitigated by averaging over multiple train/test splits, but this can get computationally expensive.

Drawback 2: only *some* of the data points are used to fit the model.

- Statistical methods tend to perform worse with less data.
- So the test-set error rate is *biased*: it tends to *overestimate* the error rate we'd see if we trained on the full data set.
- This effect is especially strong for complex models.

Leave-one-out cross validation (LOOCV)

Key idea: average over all possible testing sets of size $n_{\text{test}} = 1$.

For i in 1 to N :

- Fit the model using all data points except i .
- Predict y_i with \hat{y}_i , using the fitted model.
- Calculate the error, $\text{Err}(y_i, \hat{y}_i)$.

Estimate the error rate as:

$$\text{CV}_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{Err}(y_i, \hat{y}_i)$$

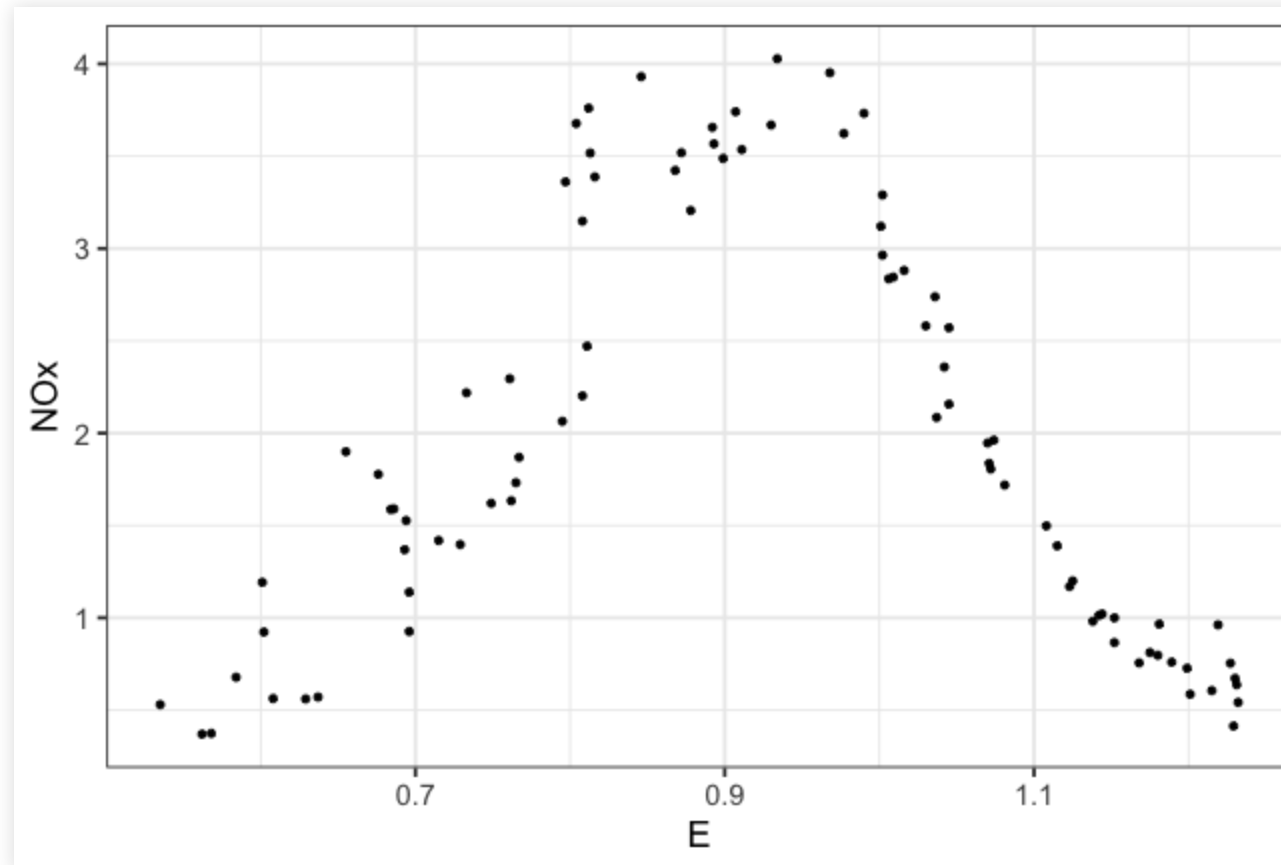
LOOCV: quick example

The data in “ethanol.csv” is from an experiment where an ethanol-based fuel was burned in a one-cylinder combustion engine.

- The experimenter varied the equivalence ratio (E), which measures the richness of the fuel-air mixture in the combustion chamber.
- For each setting of E , emission of nitrogen oxides (NO_x) were recorded.

LOOCV: quick example

```
library(tidyverse)
ethanol = read.csv('ethanol.csv')
ggplot(ethanol) + geom_point(aes(x=E, y=NOx)) +
  theme_bw(base_size=18)
```



LOOCV: quick example

KNN with $K = 5$:

```
library(FNN)

N = nrow(ethanol)
err_save = rep(0, N)
for(i in 1:N) {
  X_train = ethanol$E[-i]
  y_train = ethanol$NOx[-i]
  X_test = ethanol$E[i]
  y_test = ethanol$NOx[i]
  knn_model = knn.reg(X_train, X_test, y_train, k = 5)
  yhat_test = knn_model$pred
  err_save[i] = (y_test - yhat_test)
}
# RMSE
sqrt(mean(err_save^2))
```

```
[1] 0.3365111
```

LOOCV: quick example

Comparison with a 5th degree polynomial:

```
N = nrow(ethanol)
err_save2 = rep(0, N)
for(i in 1:N) {
  y_test = ethanol$NOx[i]
  poly5 = lm(NOx ~ poly(E, 5), data=ethanol[-i,])
  yhat_test = predict(poly5, newdata=ethanol[i,])
  err_save2[i] = (y_test - yhat_test)^2
}
# RMSE
sqrt(mean(err_save2^2))
```

```
[1] 0.2319733
```

LOOCV: pros/cons

Less bias than the train/test split approach:

- fitting with $n - 1$ points is *almost* the same as fitting with n points.
- Thus LOOCV is less prone to overestimating the training-set error.

No randomness:

- Estimating the error rate using train/test splits will yield different results when applied repeatedly (Monte Carlo variability)
- LOOCV will give the same result every time.

Downside: must re-fit n times!

K-fold cross validation

Randomly divide the data set into K nonoverlapping groups, or *folds*, of roughly equal size.

For fold $k = 1$ to K :

- Fit the model using all data points not in fold i .
- For all points (y_i, x_i) in fold k , predict \hat{y}_i using the fitted model.
- Calculate Err_k , the average error on fold k .

Estimate the error rate as:

$$\text{CV}_{(K)} = \frac{1}{K} \sum_{k=1}^K \text{Err}_k$$

10-fold CV: quick example

Back to the ethanol data:

```
N = nrow(ethanol)

# Create a vector of fold indicators
K = 10
fold_id = rep_len(1:K, N) # repeats 1:K over and over again
fold_id = sample(fold_id, replace=FALSE) # permute the order randomly

err_save = rep(0, K)
for(i in 1:K) {
  train_set = which(fold_id != i)
  y_test = ethanol$NOx[-train_set]
  poly5 = lm(NOx ~ poly(E, 5), data=ethanol[train_set,])
  yhat_test = predict(poly5, newdata=ethanol[-train_set,])
  err_save[i] = mean((y_test - yhat_test)^2)
}
# RMSE
sqrt(mean(err_save))
```

```
[1] 0.3653257
```

K-fold cross validation

- Generally requires less computation than LOOCV (K refits, versus N). If N is extremely large, LOOCV is almost certainly infeasible.
- More stable (lower variance) than running K random train/test splits.
- LOOCV is a special-case of K -fold CV (with $K = N$).

K-fold CV: bias-variance tradeoff

Key insight: there is a bias-variance tradeoff in estimating test error.

Bias comes from estimating out-of-sample error using a smaller training set than the full data set.

- LOOCV: minimal bias, since using $N - 1$ points to fit.
- K-fold: some bias, e.g. using 80% of N to fit when $K = 5$.

K-fold CV: bias-variance tradeoff

Key insight: there is a bias-variance tradeoff in estimating test error.

Variance comes from *overlap* in the training sets:

- In LOOCV, we average the outputs of N fitted models, each trained on *nearly identical* observations.
- In K-fold, we average the outputs of K fitted models that are less correlated, since they have less overlap in the training data.
Generally *less variance* than LOOCV.

Typical values: $K = 5$ to $K = 10$ (no theory; a purely empirical observation).

Back to regularization: software

There are many, many packages for fitting lasso regressions in R.

- `glmnet` is most common, and `gam1r` is my favorite. These two are very similar, and they share syntax.
- Side note: big difference is what they do beyond a simple lasso: `glmnet` does an “elastic net” penalty, while `gam1r` does a “gamma lasso” penalty. A little beyond the scope of the course, but cool stuff.
- Since we stick mostly to lasso, they’re nearly equivalent for us. `gam1r` just makes it easier to apply some model selection rules.

Both use the Matrix library representation for sparse matrices.

Sparse matrices

Often, your feature matrix will be very sparse (i.e, mostly zeros).

- This is especially true when you have lots of factors (categorical variables) as features.
- It is then efficient to ignore zero elements in actually storing X .

A simple triplet matrix is a very common storage format:

- It only stores the nonzero elements
- The row 'i', column 'j', and entry value 'v'.

Sparse matrices

For example:

$$X = \begin{bmatrix} 0 & 0 & 3.1 \\ 1.7 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

would be stored as

- $i = 1, 2$
- $j = 1, 3$
- $v = 3.1, 1.7$

Sparse matrices

For `gamlr/glmnet` you need to build your own design matrix i.e., what the $y \sim x_1 + x_2$ formula does for you inside `glm`.

We've seen how to do this with `model.matrix` for dense matrices. The sparse version `sparse.model.matrix` works the same way.

We'll walk through the example in `semiconductor.R`, which shows all these ideas:

- sparse model matrices
- regularization paths
- AIC vs. cross-validation

Dimensionality reduction

Basic idea:

- Summarize the information in the p original features into a smaller set of k *summary features* ($k < p$).
- Then try to predict y using the derived features.
- Summaries are often linear combinations of the original variables.

This is useful when:

- there are still too many variables (the haystack is too big).
- and/or the features are highly correlated and it doesn't make sense to just select a handful of them.

Dimensionality reduction: example

This poses a pretty substantial estimation problem:

- only 12 months per year, and perhaps 40 years of data.
- Yet *thousands* of other stocks in the market.

Could do variable selection from among these thousands of variables, but remember: each extra bit of hay we throw onto the haystack makes it harder to find the needles!

Dimensionality reduction: example

Suppose we want to build a model for how the returns on Apple's stock depend on all other stocks in the market. So let

- y_t = return on Apple stock in month t .
- x_{jt} = return on other stock j in month t .

A linear model would say

$$E(y_t) = \alpha + \sum_{j=1}^p \beta_j x_{jt} + e_t$$

Dimensionality reduction: example

So we try dimensionality reduction:

$$E(y_t) = \alpha + \beta m_t + e_t$$

where

- y_t is the return on a stock of interest (e.g. Apple) in period t
- m_t is the “market return”, i.e. the weighted-average return on all other stocks:

$$m_t = \sum_{j=1}^p w_j x_{jt} .$$

(The weight w_j is proportional to the size of the company.)

Dimensionality reduction: example

In finance this is called the “single index” model, a variant of the “capital asset pricing model”.

See, e.g., <https://finance.yahoo.com/quote/AAPL>

Look for their estimate of β !

Dimensionality reduction

The single-index model filter p variables into 1 variable, via a single linear combination:

$$m = \sum_{j=1}^p w_j x_j .$$

In the case of a single-index model for stocks, it is natural to take a market-weighted linear combination of all stocks. The weights come from economics, not statistics.

But what if we don't have a natural or “obvious” set of weights?

Dimensionality reduction: PCA

A very popular way to proceed here is Principal Components Analysis

- PCA is a dimensionality reduction technique that tries to represent p variables with a $k < p$ new variables.
- Like in the single-index model, these “new” variables are linear combinations of the original variables.
- The hope is that a small number of them are able to effectively represent what is going on in the original data.
- We'll study PCA in detail later. For now, we'll use it as a black box “machine” for generating summary features.

PCA: some intuition

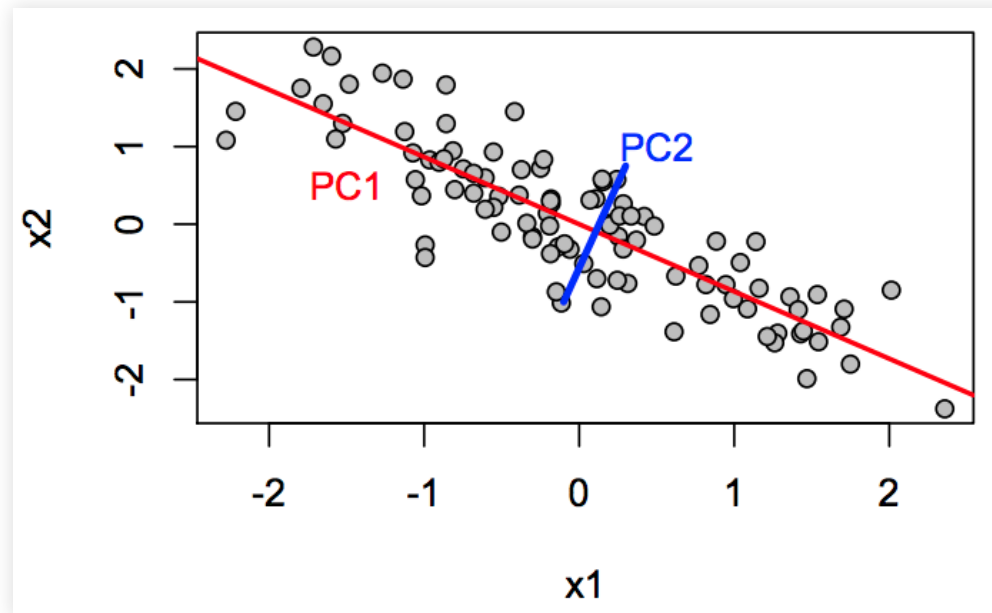
In PCA, we extract K summary variables from each feature vector x_i . Each summary $k = 1, \dots, K$ is of the form:

$$s_{ik} = \sum_{j=1}^p v_{jk} x_{ij}$$

where x_{ij} is original feature j , and s_{ik} is summary feature k , for observation i .

The coefficients v_{jk} are the “loadings” or “weights” on the original variables. The goal of PCA is to choose an “optimal” set of weights.

PCA: some intuition



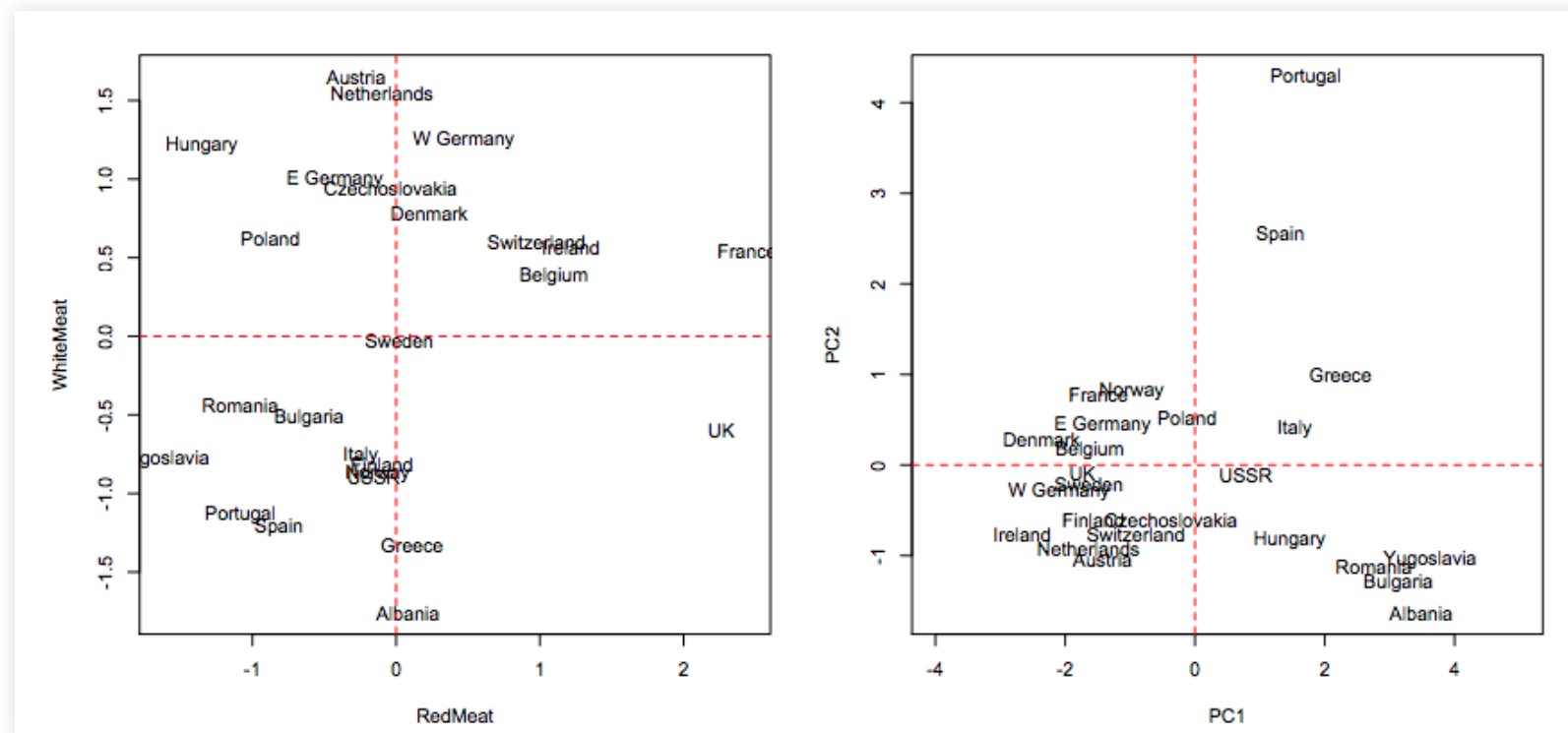
- The two variables x_1 and x_2 are very correlated. $PC1$ tells you almost everything that is going on in this dataset:

$$s_1 = 0.781x_1 - 0.625x_2$$

- PCA will look for linear combinations of the original variables that account for most of their variability.

PCA: a simple example

Data from 1970s Europe on consumption of 7 types of foods: red meat, white meat, eggs, milk, fish, cereals, starch, nuts, vegetables.



Can you interpret the summaries?

PCA: quick summary

- PCA is a great way to “collapse” many features into few.
- The choice of K (how many summaries) can be evaluated via the out-of-sample fit.
- The units of each summary variable are not interpretable in an absolute sense—only relative to each other.
- It's always a good idea to center the data before running PCA.
- Much more on PCA later!

Let's see an example in `gasoline.R`.