

Model selection and validation 1: Cross-validation

Ryan Tibshirani
Data Mining: 36-462/36-662

March 26 2013

Optional reading: ISL 2.2, 5.1, ESL 7.4, 7.10

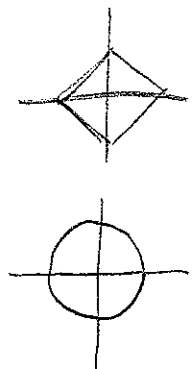
Reminder: modern regression techniques

Over the last two lectures we've investigated modern regression techniques. We saw that linear regression has generally low bias (zero bias, when the true model is linear) but high variance, leading to poor predictions. These modern methods introduce some bias but significantly reduce the variance, leading to better predictive accuracy

Given a response $y \in \mathbb{R}^n$ and predictors $X \in \mathbb{R}^{n \times p}$, we can think of these modern methods as constrained least squares:

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \quad \text{subject to } \|\beta\|_q \leq t,$$

- ▶ $q = 1$ gives the lasso (and $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$)
- ▶ $q = 2$ gives ridge regression (and $\|\beta\|_2 = \sqrt{\sum_{j=1}^p \beta_j^2}$)



Regularization at large

Most other modern methods for regression can be expressed as

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \quad \text{subject to } R(\beta) \leq t,$$

or equivalently,

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \cdot R(\beta).$$

The term R is called a penalty or regularizer, and modifying the regression problem in this way is called applying regularization

- ▶ Regularization can be applied beyond regression: e.g., it can be applied to classification, clustering, principal component analysis
- ▶ Regularization goes beyond sparsity: e.g., design R to induce smoothness or structure, instead of pure sparsity

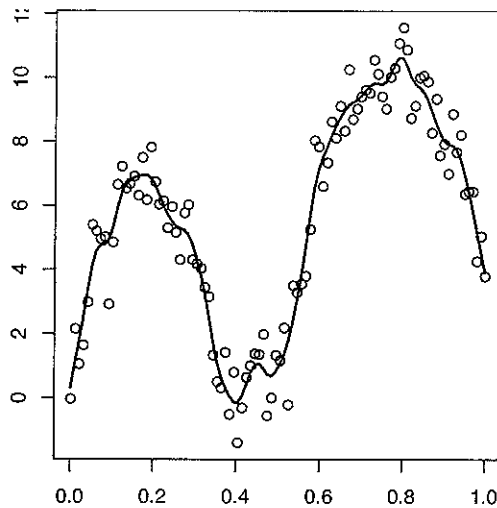
Example: smoothing splines

Smoothing splines use a form of regularization:

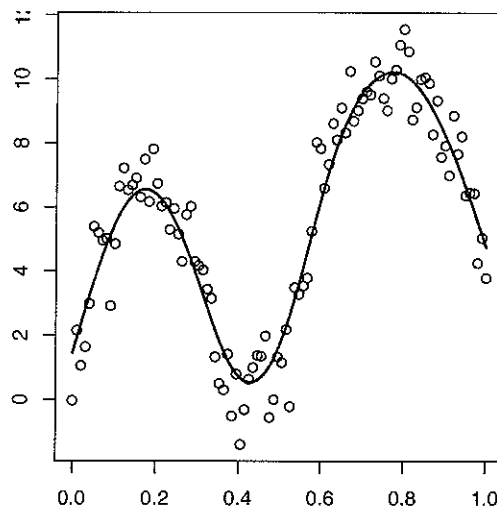
Regularization

$$\hat{f} = \underset{f}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n (y_i - f(x_i))^2}_{\text{Loss}} + \lambda \cdot \underbrace{\int (f''(x))^2 dx}_{R(f) \text{ Penalty}}$$

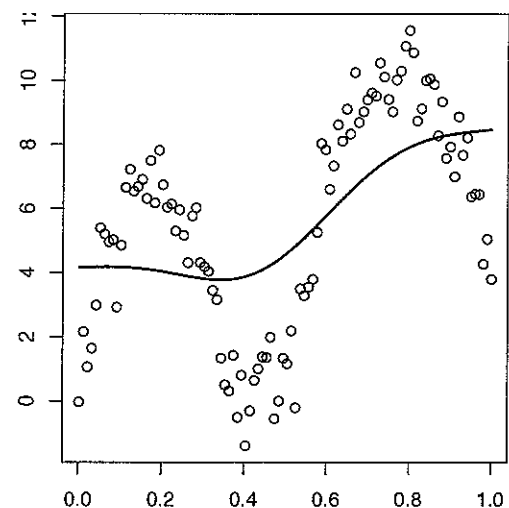
Example with $n = 100$ points:



λ too small



λ just right



λ too big

Choosing a value of the tuning parameter

Each regularization method has an associated tuning parameter: e.g., this was λ in the smoothing spline problem, and λ for lasso and ridge regression in the penalized forms (or t in the constrained forms)

The tuning parameter controls the amount of regularization, so choosing a good value of the tuning parameter is crucial. Because each tuning parameter value corresponds to a fitted model, we also refer to this task as model selection

What we might consider a good choice of tuning parameter, however, depends on whether our goal is prediction accuracy or recovering the right model for interpretation purposes. We'll cover choosing the tuning parameter for the purposes of prediction; choosing the tuning parameter for the latter purpose is a harder problem

Prediction error and test error

Our usual setup: we observe

$$y_i = f(x_i) + \epsilon_i, \quad i = 1, \dots, n$$

where $x_i \in \mathbb{R}^p$ are fixed (nonrandom) predictor measurements, f is the true function we are trying to predict (think $f(x_i) = \underline{x_i^T \beta^*}$ for a linear model), and ϵ_i are random errors

We call (x_i, y_i) , $i = 1, \dots, n$ the training data. Given an estimator \hat{f} built on the training data, consider the average prediction error over all inputs

$$\text{PE}(\hat{f}) = \text{E} \left[\frac{1}{n} \sum_{i=1}^n (y'_i - \hat{f}(x_i))^2 \right]$$

where $y'_i = f(x_i) + \epsilon'_i$, $i = 1, \dots, n$ are another set of observations, independent of y_1, \dots, y_n

Suppose that $\hat{f} = \hat{f}_\theta$ depends on a tuning parameter θ , and we want to choose θ to minimize $\text{PE}(\hat{f}_\theta)$

If we actually had training data y_1, \dots, y_n and test data y'_1, \dots, y'_n (meaning that we don't use this to build \hat{f}_θ) in practice, then we could simply use

$$\text{TestErr}(\hat{f}_\theta) = \frac{1}{n} \sum_{i=1}^n (y'_i - \hat{f}_\theta(x_i))^2$$

called the test error, as an estimate for $\text{PE}(\hat{f}_\theta)$. The larger that n is, the better this estimate

We usually don't have test data. So what to do instead?

What's wrong with training error?

It may seem like

$$\text{TestErr}(\hat{f}_\theta) = \frac{1}{n} \sum_{i=1}^n (y'_i - \hat{f}_\theta(x_i))^2 \quad \text{and}$$

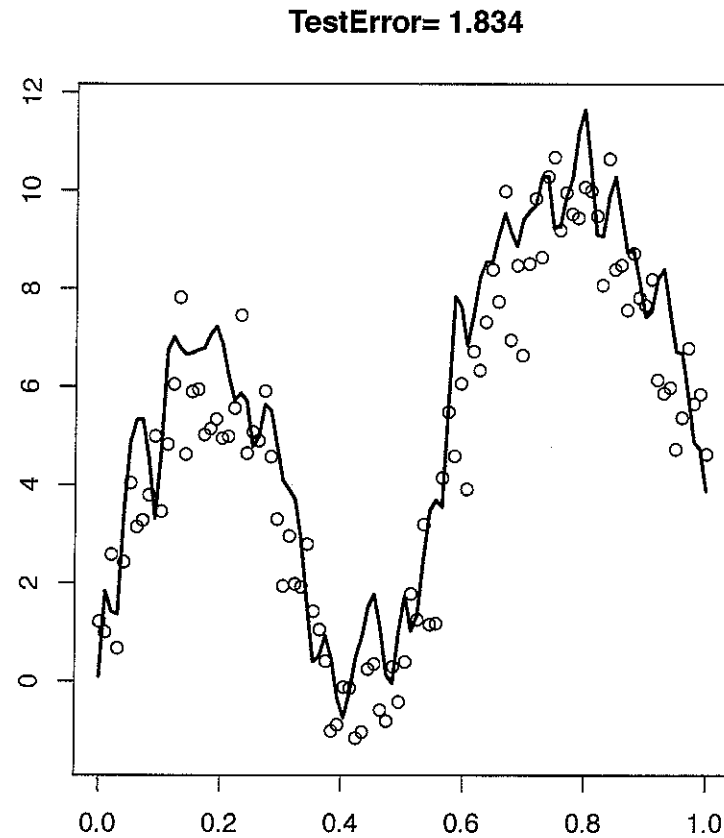
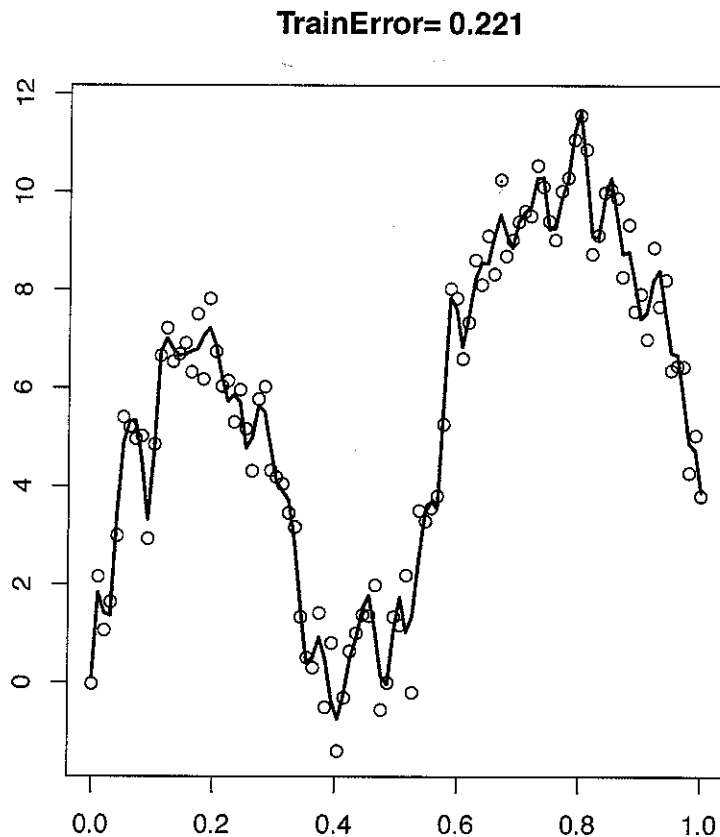
$$\text{TrainErr}(\hat{f}_\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_\theta(x_i))^2$$

shouldn't be too different—after all, y_i and y'_i are independent copies of each other. The second quantity is called the training error: this is the error of \hat{f} as measured by the data we used to build it

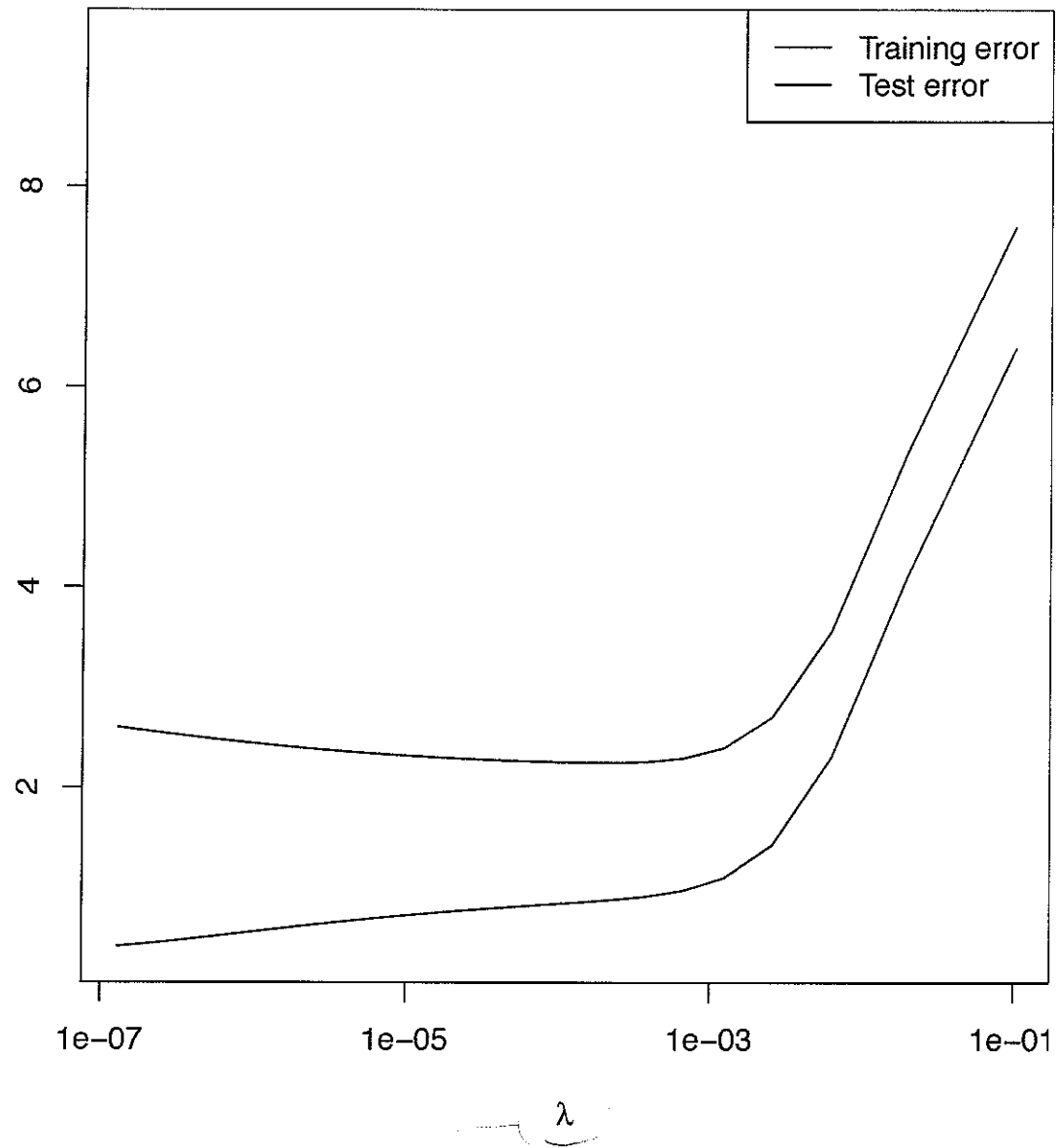
But actually, the training and test error curves are fundamentally different. Why?

Example: smoothing splines (continued)

Back to our smoothing spline example: for a small value of the tuning parameter λ , the training and test errors are very different



Training and test error curves, averaged over 100 simulations:

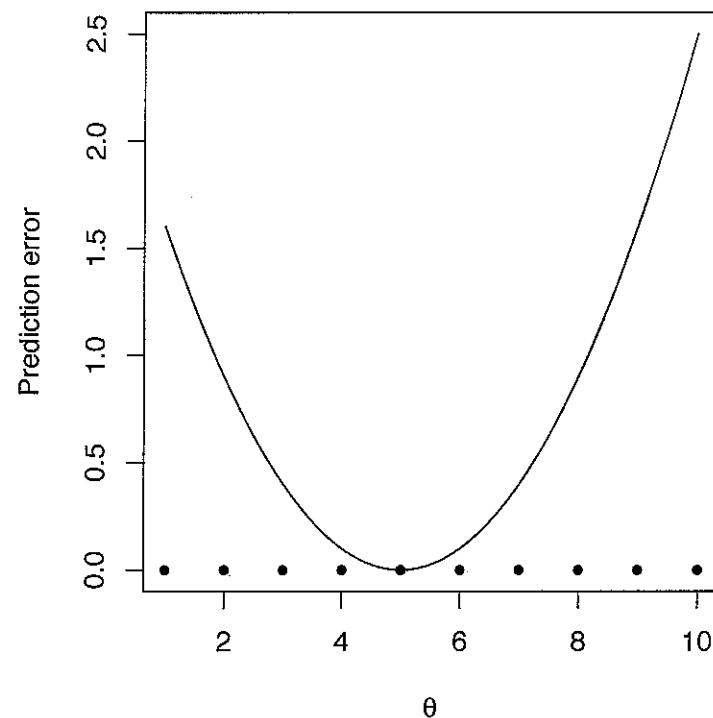


Cross-validation

Cross-validation is a simple, intuitive way to estimate prediction error

Given training data (x_i, y_i) , $i = 1, \dots, n$ and an estimator \hat{f}_θ , depending on a tuning parameter θ

Even if θ is a continuous parameter, it's usually not practically feasible to consider all possible values of θ , so we discretize the range and consider choosing θ over some discrete set $\{\theta_1, \dots, \theta_m\}$



For a number K , we split the training pairs into K parts or “folds”
(commonly $K = 5$ or $K = 10$)

10-fold CV

| 1 | 2 | 3 | 4 | 5 |
|------------------------------|--------------|-------------------|--------------|--------------|
| Train (x_i, y_i) | Train | Validation | Train | Train |

K -fold cross validation considers training on all but the k th part, and then validating on the k th part, iterating over $k = 1, \dots, K$

(When $K = n$, we call this leave-one-out cross-validation, because we leave out one data point at a time)

K -fold cross validation procedure:

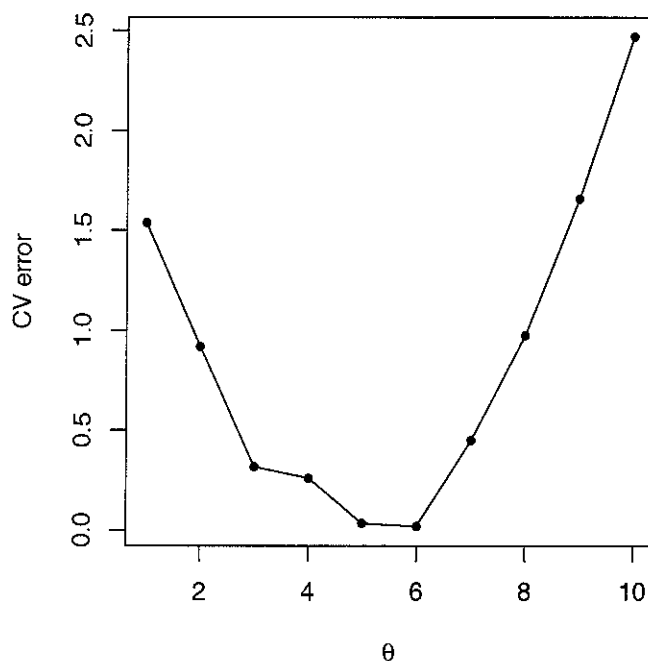
- ▶ Divide the set $\{1, \dots, n\}$ into K subsets (i.e., folds) of roughly equal size, F_1, \dots, F_K randomly.
- ▶ For $k = 1, \dots, K$:
 - ▶ Consider training on (x_i, y_i) , $i \notin F_k$, and validating on (x_i, y_i) , $i \in F_k$
 - ▶ For each value of the tuning parameter $\theta \in \{\theta_1, \dots, \theta_m\}$, compute the estimate \hat{f}_{θ}^{-k} on the training set, and record the total error on the validation set:

$$\underline{e_k(\theta)} = \sum_{i \in F_k} (y_i - \hat{f}_{\theta}^{-k}(x_i))^2$$

- ▶ For each tuning parameter value θ , compute the average error over all folds,

$$\text{CV}(\theta) = \frac{1}{n} \sum_{k=1}^K e_k(\theta) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in F_k} (y_i - \hat{f}_{\theta}^{-k}(x_i))^2$$

Having done this, we get a cross-validation error curve $CV(\theta)$ (this curve is a function of θ), e.g.,



and we choose the value of tuning parameter that minimizes this curve,

$$\hat{\theta} = \underset{\theta \in \{\theta_1, \dots, \theta_m\}}{\operatorname{argmin}} CV(\theta)$$

set.seed

Example: choosing λ for the lasso

Recall our running example from last time: $n = 50$, $p = 30$, and the true model is linear with 10 nonzero coefficients. We consider

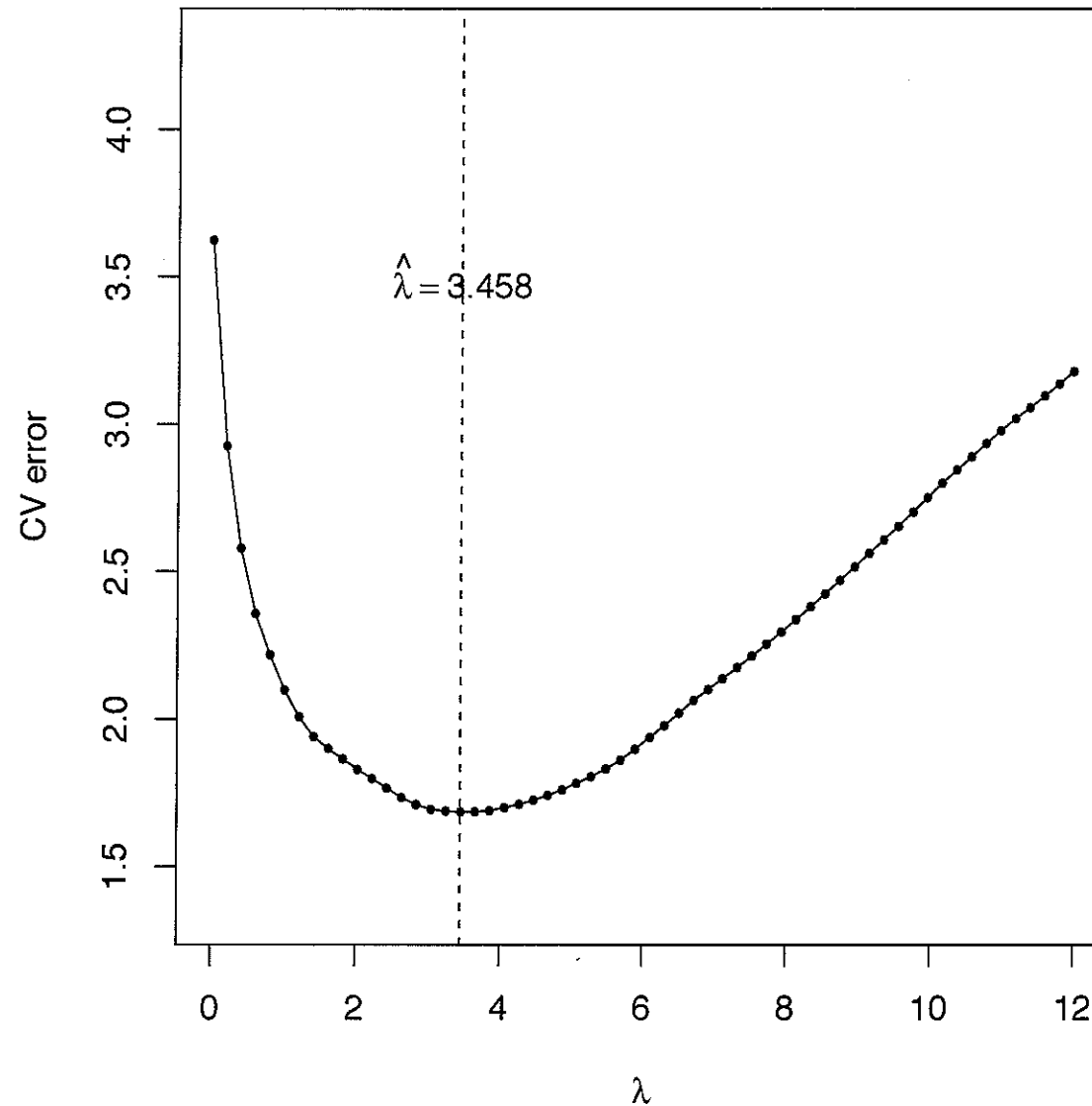
$$\hat{\beta}^{\text{lasso}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

and use 5-fold cross-validation to choose λ . Sample code for making the folds:

```
K = 5 set.seed(0)
i.mix = sample(1:n)
folds = vector(mode="list",length=K)
# divide i.mix into 5 chunks of size n/K=10, and
# store each chunk in an element of the folds list }
> folds
[[1]]
 [1] 30 29  4  2 50 42 27 25 23 41

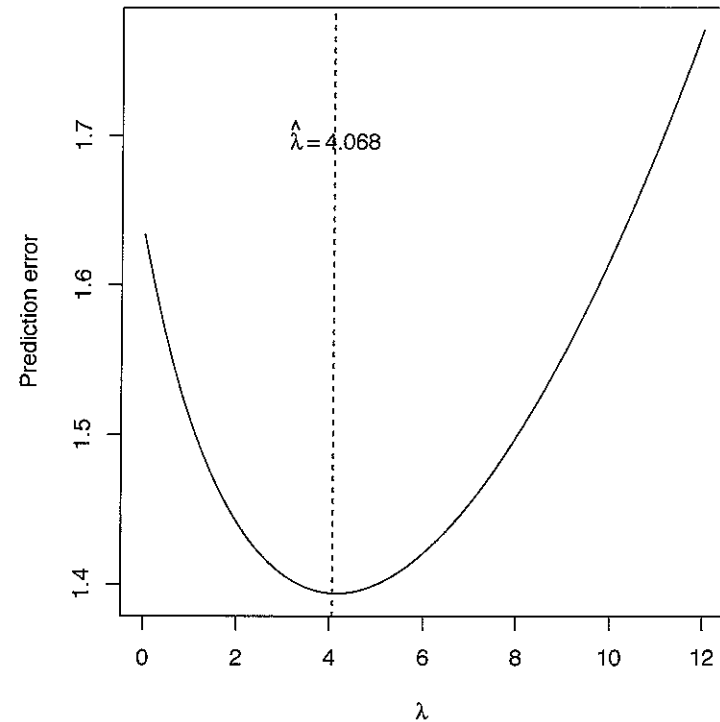
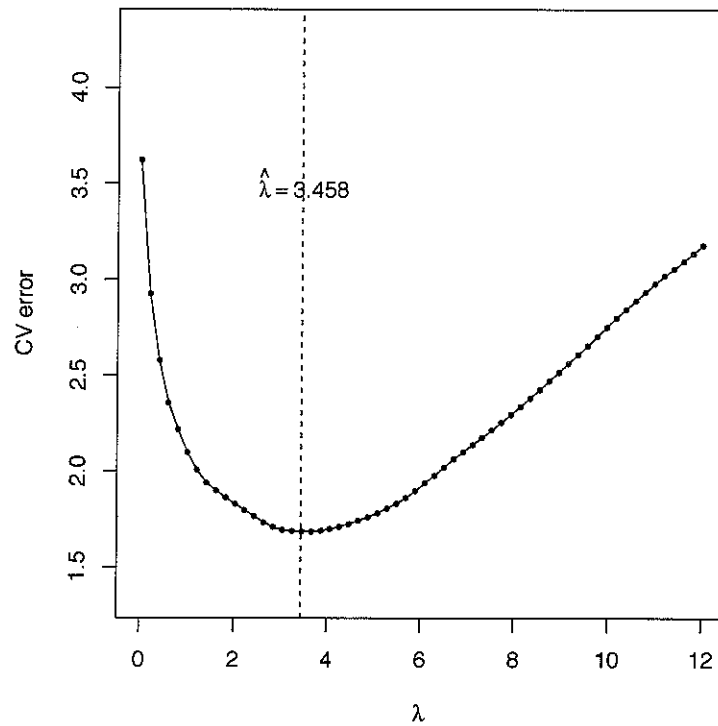
[[2]]
 [1] 21 44 45  9 10 26 16  6 24 48
...
```

We consider the values: $\lambda_{\text{am}} = \text{seq}(0, 12, \text{length}=60)$. The resulting cross-validation error curve:



$\hat{\beta}^{\text{lasso}}$ at $\lambda = 3.458$

How does this compare to the prediction error curve? Because this is a simulation, we can answer this question.



Standard errors for cross-validation

One nice thing about K -fold cross-validation (for a small $K \ll n$, e.g., $K = 5$) is that we can estimate the standard deviation of $CV(\theta)$, at each $\theta \in \{\theta_1, \dots, \theta_m\}$

First, we just average the validation errors in each fold:

$$CV_k(\theta) = \frac{1}{n_k} e_k(\theta) = \frac{1}{n_k} \sum_{i \in F_k} (y_i - \hat{f}_{\theta}^{-k}(x_i))^2$$

where n_k is the number of points in the k th fold

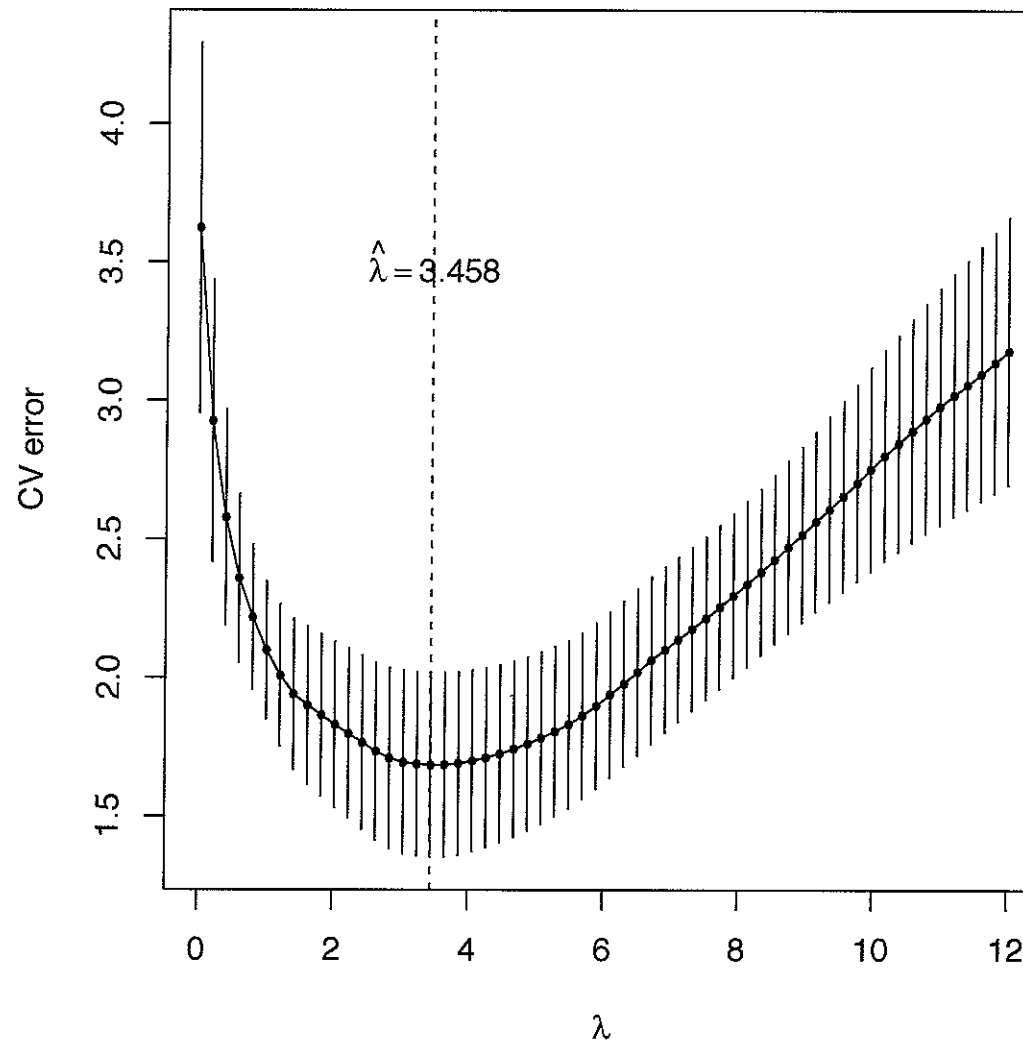
We then compute the sample standard deviation of $CV_1(\theta), \dots, CV_K(\theta)$,

$$SD(\theta) = \sqrt{\text{var}(CV_1(\theta), \dots, CV_K(\theta))} \quad \text{sd}(\)$$

Finally we estimate the standard deviation of $CV(\theta)$ —called the standard error of $CV(\theta)$ —by $SE(\theta) = SD(\theta)/\sqrt{K}$

Example: choosing λ for the lasso (continued)

The cross-validation error curve from our lasso example, with \pm standard errors:



The one standard error rule

The one standard error rule is an alternative rule for choosing the value of the tuning parameter, as opposed to the usual rule

$$\hat{\theta} = \underset{\theta \in \{\theta_1, \dots, \theta_m\}}{\operatorname{argmin}} \operatorname{CV}(\theta)$$

We first find the usual minimizer $\hat{\theta}$ as above, and then move θ in the direction of increasing regularization (this may be increasing or decreasing θ , depending on the parametrization) as much as we can, such that the cross-validation error curve is still within one standard error of $\operatorname{CV}(\hat{\theta})$. In other words, we maintain

$$\operatorname{CV}(\theta) \leq \operatorname{CV}(\hat{\theta}) + \operatorname{SE}(\hat{\theta})$$

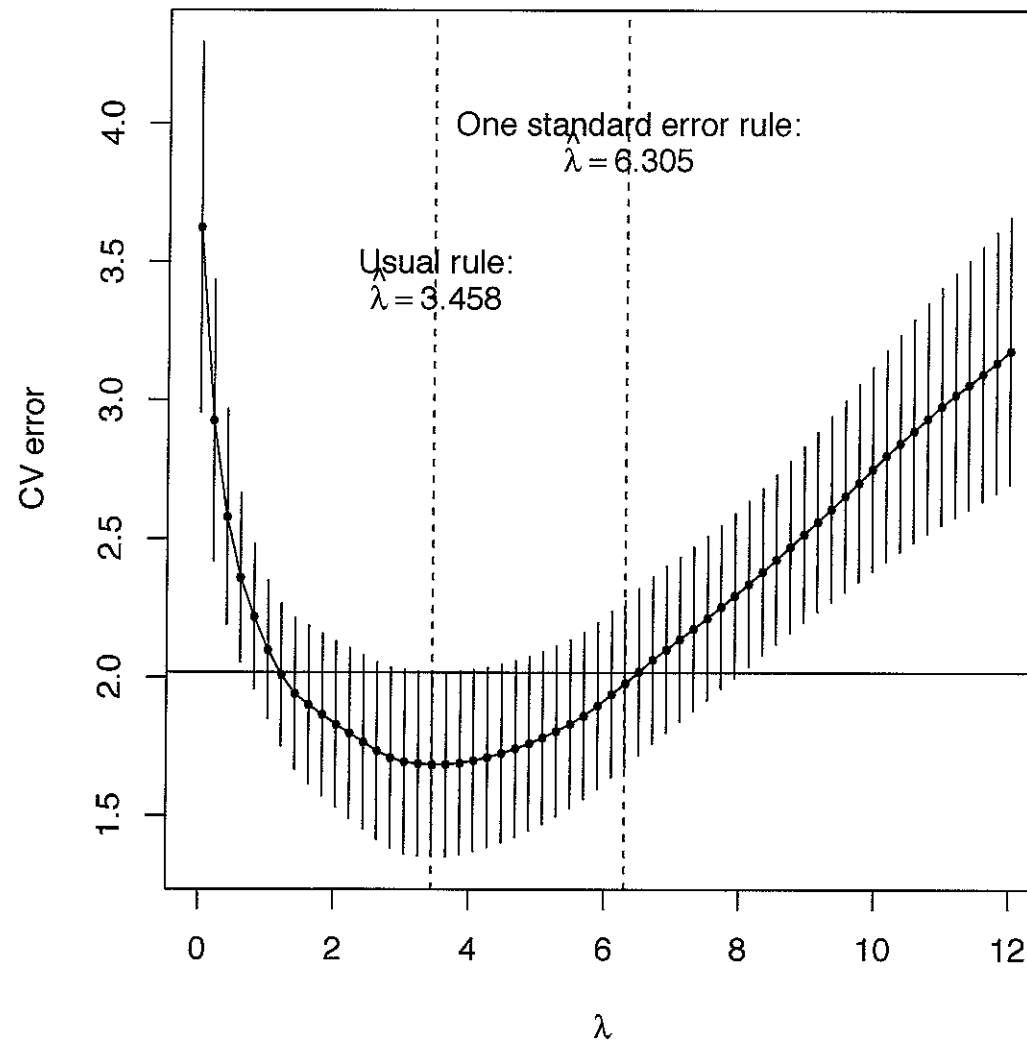
Idea: “All else equal (up to one standard error), go for the simpler (more regularized) model”

parametrize

$$\|\beta\|_1 \leq t$$

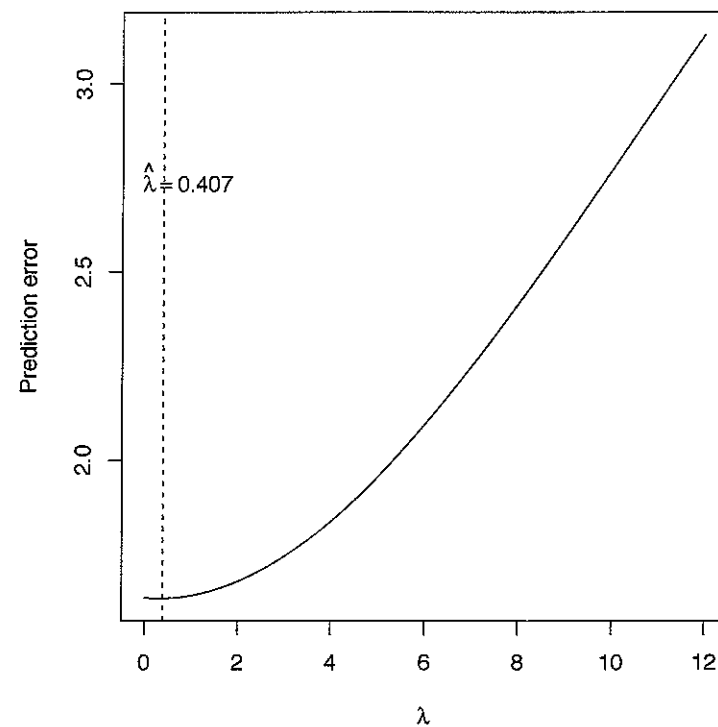
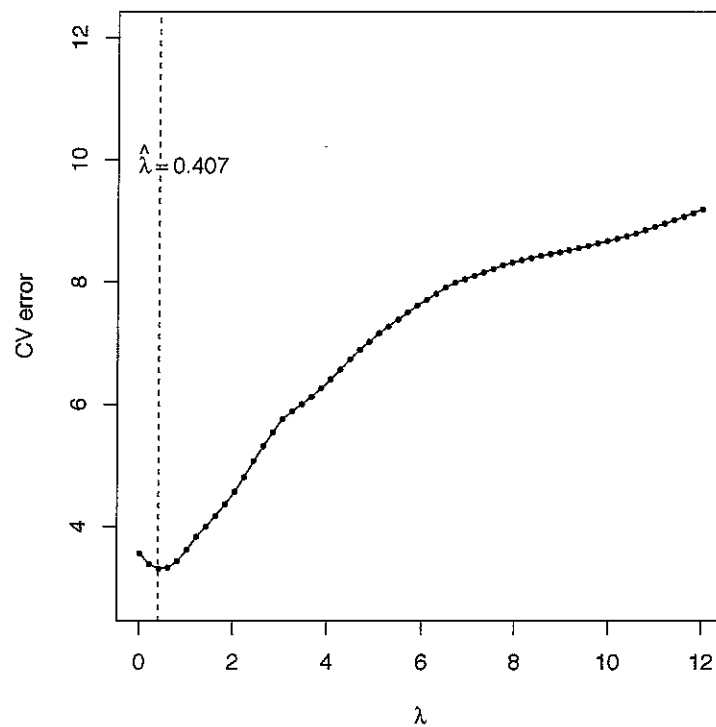
Example: choosing λ for the lasso (continued)

Back to our lasso example: the one standard error rule chooses a larger value of λ (more regularization):

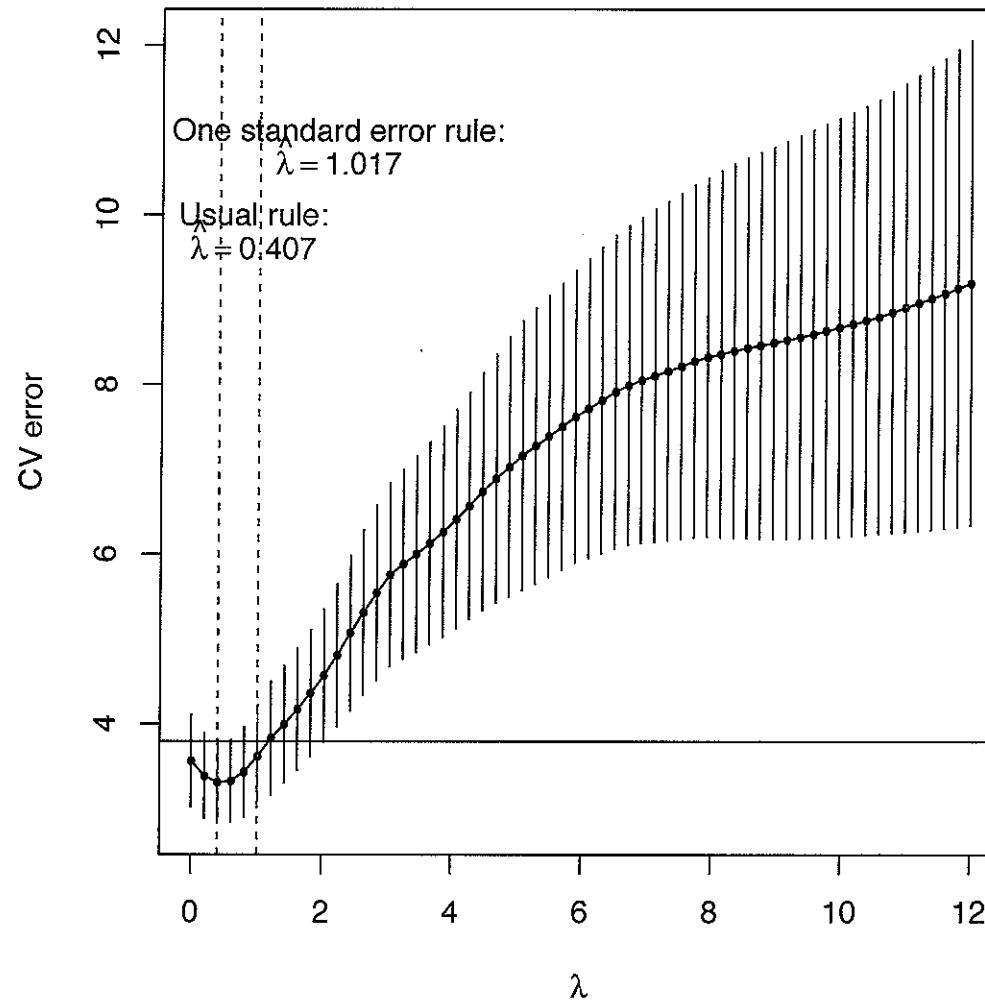


What happens if we really shouldn't be shrinking in the first place?
We'd like cross-validation, our automated tuning parameter selection procedure, to choose a small value of λ

Recall our other example from last time: $n = 50$, $p = 30$, and the true model is linear with all moderately large coefficients



The standard errors here reflect uncertainty about the shape of the CV curve for large λ



The one standard error rule doesn't move much from the minimum

Prediction versus interpretation

Remember: cross-validation error estimates prediction error at any fixed value of the tuning parameter, and so by using it, we are implicitly assuming that achieving minimal prediction error is our goal

Choosing λ for the goal of recovering the true model, for the sake of interpretation, is somewhat of a different task. The value of the parameter that achieves the smallest cross-validation error often corresponds to not enough regularization for these purposes. But the one standard error rule is a step in the right direction

There are other (often more complicated) schemes for selecting a value of the tuning parameter for the purposes of true model recovery

Recap: cross-validation

Training error, the error of an estimator as measured by the data used to fit it, is not a good surrogate for prediction error. It just keeps decreasing with increasing model complexity

Cross-validation, on the other hand, much more accurately reflects prediction error. If we want to choose a value for the tuning parameter of a generic estimator (and minimizing prediction error is our goal), then cross-validation is the standard tool

We usually pick the tuning parameter that minimizes the cross-validation error curve, but we can also employ the one standard error rule, which helps us pick a more regularized model

Next time: model assessment and more cross-validation

How would we do cross-validation for a structured problem like that solved by smoothing splines?

