**Douglas Nychka,** www.image.ucar.edu/~nychka

- Locally weighted averages
- Penalized least squares smoothers
- Properties of smoothers
- Splines and Reproducing Kernels
- The interpolation proof
- CV and the smoothing parameter





## The additive statistical model:

Given n pairs of observations  $(x_i, y_i)$ , i = 1, ..., n

$$y_i = g(x_i) + \epsilon_i$$

 $\epsilon_i$ 's are random errors and g is an unknown smooth function.

The goal is to estimate a function g based on the observations

## A 2-d example

Predict surface ozone where it is not monitored.



## Local linear surface estimates

Use local information to predict unobserved values

Use a linear regression based on close by observations.

$$y_i = \beta_1 + \mathbf{lon}_i \beta_2 + \mathbf{lat}_i \beta_3 + \epsilon_i$$

find  $\hat{\beta}$  by least squares.

The prediction at location (-88, 41) is just a weighted average of the observations.

$$\hat{g}(-88,41) = \hat{\beta}_1 + -88\hat{\beta}_2 + 41\hat{\beta}_3 = \sum_{i=1,n} w_i y_i$$

## A kernel estimator

Determine the weights based on the distance to the prediction points

$$w_i \sim (1/h) K((x-x_i)/h)$$

(and normalize so that the weights sum to one.)

Kernel: KK is bump shaped e.g. a normal

#### Bandwidth: h

h controls the spread of K as h gets large the estimate is just the average.

## Some kernel estimates for ozone



## Linear smoothers

Let  $\hat{g} = g(x_1), ..., g(x_n)$  be the prediction vector at the observed points.

# A smoother matrix satisfies $\hat{g} = Ay$ where

- A is an  $n \times n$  matrix
- eigenvalues of A are in the range [0,1].

## **Note:** $||Ay|| \le ||y||$

Usually values in between the data are filled in by interpolating the predictions at the observations.

## Problems with local regression and kernels

- How large should the neighborhood/bandwidth be?
- What is the uncertainty of the prediction?
- Predicting in between observations is *ad hoc* and can get weird when the error is small.

#### Problems with local regression and kernels

- How large should the neighborhood/bandwidth be?
- What is the uncertainty of the prediction?
- Predicting in between observations is *ad hoc* and can get weird when the error is small.

But the theoretical properties of kernel estimators are well understood ...

$$E\left[g(\boldsymbol{x}) - \hat{g}(\boldsymbol{x})\right]^2 = \boldsymbol{h}^4 K_2 / 4 + \frac{\sigma^2}{n\boldsymbol{h}} K_0$$

$$MSE = Bias^2 + Variance$$

#### **Ridge regression**

Start with your favorite n basis functions  $\{\psi_k\}_{k=1}^n$  The estimate has the form

$$f(x) = \sum_{l=1}^{n} \theta_k \psi_k(x)$$

where  $\theta = (\theta_1, \ldots, \theta_n)$  are the coefficients.

Let  $W_{i,k} = \psi_k(x_i)$  so  $f = W\theta$ 

Estimate the coefficients by a penalized least squares

## Sum of squares( $\theta$ ) + penalty on $\theta$ $\min_{\theta} \sum_{i=1}^{n} (y - [W\theta]_i)^2 + \lambda \theta^T B \theta$

with  $\lambda > 0$  a hyperparameter and B a nonnegative definite matrix.

Estimate the coefficients by a penalized least squares

## Sum of squares( $\theta$ ) + penalty on $\theta$ $\min_{\theta} \sum_{i=1}^{n} (y - [W\theta]_i)^2 + \lambda \theta^T B \theta$

with  $\lambda > 0$  a hyperparameter and B a nonnegative definite matrix.

or in general,

- log likelihood +  $\lambda$  penalty on  $\theta$ 

In any case once we have the parameter estimates these can be used to evaluate  $\hat{g}$  at any point.

Just calculus ...

- Take derivatives of the penalized likelihood w/r to  $\theta$ ,
- set equal to zero,
- $\bullet$  solve for  $\theta$

Just calculus ...

- Take derivatives of the penalized likelihood w/r to  $\theta$ ,
- set equal to zero,
- solve for  $\theta$

The monster ...

$$\hat{\boldsymbol{\theta}} = (W^T W + \lambda B)^{-1} W^T \boldsymbol{y}$$

 $\hat{g} = W\hat{\theta} = W(W^TW + \lambda B)^{-1}W^Ty = A(\lambda)y$ 

Why is this a smoothing matrix?

If W is symmetric

 $A(\lambda) = W(W^{T}W + \lambda B)^{-1}W^{T} = (I + \lambda W^{-1}BW^{-1})^{-1}$ 

## Effective degrees of freedom in the smoother

For linear regression trace  $A(\lambda)$  gives us the number of parameters. (Because it is a projection matrix)

By analogy,  $trA(\lambda)$  is measure of the effective number of degrees of freedom attributed to the smooth surface For linear regression trace  $A(\lambda)$  gives us the number of parameters. (Because it is a projection matrix)

By analogy,  $trA(\lambda)$  is measure of the effective number of degrees of freedom attributed to the smooth surface

## A useful decomposition Recall: $A(\lambda) = (I + \lambda W^{-1} B W^{-1})^{-1}$

One can always find an orthogonal matrix, U so that  $U^T U = I$  and  $(W^{-1}BW^{-1}) = U\Gamma U^T$ 

where  $\[Gamma]$  is diagonal.

$$A(\lambda) = (I + \lambda U \Gamma U^T)^{-1} = U(I + \lambda \Gamma)^{-1} U^T$$

A simple formula for the trace

So

$$\mathbf{tr}A(\lambda) = \mathbf{tr}U(I + \lambda\Gamma)^{-1}U^T = \mathbf{tr}(I + \lambda\Gamma)^{-1}U^T U$$
$$= \mathbf{tr}(I + \lambda\Gamma)^{-1} = \sum_{i=1}^n \frac{1}{1 + \lambda\Gamma_{ii}}$$

The Gamma's are all nonnegative so this must be increasing as  $\lambda$  decreases.

The relationship is one-to-one with  $\lambda$  and independent of the data so we can always talk about  $\lambda$  in terms of the effective degrees of freedom.

## **Splines**

One obtains a spline estimate using a specific basis and a specific penalty matrix. Splines are confusing because the basis is a bit mysterious.

The classic cubic smoothing spline: For curve smoothing in one dimension,

$$\min_{f} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \int (f''(x))^2 dx$$

The second derivative measures the roughness of the fitted curve.

## **S**plines

One obtains a spline estimate using a specific basis and a specific penalty matrix. Splines are confusing because the basis is a bit mysterious.

The classic cubic smoothing spline: For curve smoothing in one dimension,

$$\min_{f} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \int (f''(x))^2 dx$$

The second derivative measures the roughness of the fitted curve.

The solution, is continuous up to its second derivative and is a piecewise cubic polynomial in between the observation points.

Where does this come from?

## **Climate for Colorado**



## Cubic splines with different $\lambda$ s



## Some abstraction

## **Reproducing kernels**

Think of these as covariance functions ... k(x, x')To get basis functions we will hold one argument fixed at some value of x and let the other vary.

## A Space of functions

Let  $\mathcal{H}$  be the smallest Hilbert space with inner product  $\langle \cdot, \cdot \rangle$  such that  $k(x, \cdot) \in \mathcal{H}$  for all x and

$$\langle k(x,\cdot), k(x',\cdot) \rangle = k(x,x')$$

k reproduces itself!

## Some abstraction

## **Reproducing kernels**

Think of these as covariance functions ... k(x, x')To get basis functions we will hold one argument fixed at some value of x and let the other vary.

## A Space of functions

Let  $\mathcal{H}$  be the smallest Hilbert space with inner product  $\langle \cdot, \cdot \rangle$  such that  $k(x, \cdot) \in \mathcal{H}$  for all x and

$$\langle k(x,\cdot), k(x',\cdot) \rangle = k(x,x')$$

k reproduces itself!

For 
$$f$$
 in  $\mathcal{H}$ ,  $\langle k(x, \cdot), f \rangle = f(x)$  !

# A variational problem With $\mathcal{H}$ and its inner product,

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n (y_i - f(\boldsymbol{x}_i))^2 + \lambda \langle f, f \rangle$$

# A variational problem With $\mathcal{H}$ and its inner product,

$$\min_{f \in \mathcal{H}} \sum_{i=1}^{n} (y_i - f(\boldsymbol{x}_i))^2 + \lambda \langle f, f \rangle$$

#### The solution

Choose basis functions  $k(., x_i)$  for  $1 \le i \le n$ 

$$\widehat{f}(x) = \sum_{i=1}^{n} \widehat{\theta}_i k(x, x_i)$$

Choose roughness penalty matrix B = W

# A variational problem With $\mathcal{H}$ and its inner product,

$$\min_{f \in \mathcal{H}} \sum_{i=1}^{n} (y_i - f(\boldsymbol{x}_i))^2 + \lambda \langle f, f \rangle$$

#### The solution

Choose basis functions  $k(., x_i)$  for  $1 \le i \le n$ 

$$\widehat{f}(x) = \sum_{i=1}^{n} \widehat{\theta}_i k(x, x_i)$$

Choose roughness penalty matrix B = Wthis gives

$$\hat{\boldsymbol{\theta}} = (WW^T + \lambda W)^{-1} W^T \boldsymbol{y} = (W + \lambda I)^{-1} \boldsymbol{y}$$

$$\widehat{g} = W\widehat{\theta} = W(W + \lambda I)^{-1}y = (I + \lambda W^{-1})^{-1}y$$

## Covariances

Here  $W_{ij} = k(x_i, x_j)$ this looks like a covariance matrix doesn't it.

## The Proof

The main part is to use the minimum interpolation properties of splines. It is easiest to prove this in more generality for any reproducing kernel Hilbert space.

A simple proof is to guess at the minimizing solution and use the optimal interpolation results to show that there can not be another solution. – more on this later.

## A 1-d cubic smoothing spline

The key step is a decomposition:  $f(x) = \beta_1 + \beta_2 x + h(x)$ 

$$\min_{f} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \int (f''(x))^2 dx$$

becomes

$$\min_{\beta,h} \sum_{i=1}^{n} (y_i - \beta_1 + \beta_2 x_i + h(x_i))^2 + \lambda \int (h''(x))^2 dx$$

 $\mathcal{H}$ :

#### All functions with 2 derivatives

Inner product: related to the integral

Usual (Wahba style) Reproducing Kernel:  $k(x, x') = |x - x'|^3 + \text{linear terms}$ 

# A less scary reproducing kernel for the same cubic spline problem

On the interval [0, 1]

Let 
$$k(u, v) = u^2 v/2 - u^3/6$$
 for  $u < v$   
and

 $k(u,v) = v^2 u/2 - v^3/6$  for  $u \ge v$ 

k(.,v) for fixed v, is a function that is a cubic polynomial from 0 to v and is then a linear function for u > v. It is twice differentiable.

Let  $\mathcal{H}$  be a Hilbert space of functions on [0,1] where f(0) = 0 and f'(0) = 0 and with inner product

$$\langle f,g\rangle = \int_0^1 f''(u)g''(u)du$$

k is a reproducing kernel for this space! Show this just using integration by parts.

### More on cubic splines

To be explicit the cubic smoothing spline solution has the form:

$$\widehat{f}(x) = \widehat{\beta}_1 + \widehat{\beta}_2 x + \sum_{i=1}^n \widehat{\theta}_i k(x, x_i)$$

Besides verifying that the reproducing on the previous page works with the integrated second derivative inner product one can also check that it is the covariance function for integrated Brownian motion,B(t):

$$X(u) = \int_0^u B(t) dt,$$

with B(0)=0. i.e.

$$E(X(u)X(v)) = k(u,v)$$

So as we will see later, the Bayesian prior associated with the cubic smoothing spline is an integrated Brownian motion.

## The last details ...

Given the reproducing kernel, the problem is

$$\min_{\boldsymbol{\beta},\boldsymbol{\theta}} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{W}\boldsymbol{\theta}||^2 + \lambda \boldsymbol{\theta}^T \boldsymbol{W}\boldsymbol{\theta}$$

*X* is the regression matrix with columns 1 and  $\{x_i\}$ First minimize over  $\theta$  with  $\beta$  fixed. Plug in solution and then minimize over  $\beta$ .

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

 $\Omega = (W + \lambda I)$ 

$$\hat{\theta} = (W + \lambda I)^{-1} (y - X\hat{\beta})$$

## A 2-d thin plate smoothing spline

$$\min_{f} \sum_{i=1}^{n} (y_i - f_i)^2 + \lambda \int_{\Re^2} \left(\frac{\partial^2 f}{\partial^2 u}\right)^2 + 2\left(\frac{\partial^2 f}{\partial u \partial v}\right)^2 + \left(\frac{\partial^2 f}{\partial^2 v}\right)^2 du dv$$

Collection of second partials is invariant to a rotation.

Again, separate off the linear part of f.  $f(x) = \beta_1 + \beta_2 x_1 + \beta_3 x_2 + h(x)$ 

## **Reproducing Kernel:**

 $k(x, x') = ||x - x'||^2 log(||x - x'||) + linear terms$ 

leading to basis functions that are bumps at the observation locations.

## Some thin plate splines for the ozone data



## The interpolation problem: a editorial

## Splines arose in numerical analysis with the interpolation problem:

Find a curve g that minimizes  $\int (g''(x))^2 dx$ , subject to  $g(x_i) = y_i$ 

## The interpolation problem: a editorial

## Splines arose in numerical analysis with the interpolation problem:

Find a curve g that minimizes  $\int (g''(x))^2 dx$ , subject to  $g(x_i) = y_i$ 

We already know how to do this by letting  $\lambda \to 0$ 

- Statisticians are, of course, suspicious of fitting data exactly
- 1-d splines have some very fast computational properties – that do not extend to higher dimensions!
- Assuming a reproducing kernel is equivalent to a model for the unknown function. What is that model?

## The general interpolation problem

Given  $\{x_i, y_i\}$  and a reproducing kernel Hilbert space. Find a function g in the space so that minimizes  $\langle g, g \rangle$ . subject to the constraints  $g(x_i) = y_i$ ,  $1 \le i \le n$ 

The solution

$$g = \sum_{i=1}^{n} \theta_i k(., x_i)$$

Where  $\theta$  is determined by solving a system of n linear equations to guarentee the interpolation constraints.

## **Interpolation Proof**

#### Suppose

t

$$g = \sum_{i=1}^{n} \theta_i k(., x_i)$$

interpolates  $\{x_i, y_i\}$  and so does some other function h. We will show that  $\langle h, h \rangle > \langle g, g \rangle$ . Let  $\delta = h - a$  and so  $\langle \delta, \delta \rangle > 0$ 

$$\langle h,h\rangle = \langle g+\delta,g+\delta\rangle = \langle g,g\rangle + 2\langle g,\delta\rangle + \langle \delta,\delta\rangle$$

The cool part, using the reproducing property

$$\langle g, \delta \rangle = \langle \sum_{i=1}^{n} \theta_{i} k(., x_{i}), \delta \rangle = \sum_{i=1}^{n} \theta_{i} \langle k(., x_{i}), \delta \rangle = \sum_{i=1}^{n} \theta_{i} \delta(x_{i}) = 0$$
  
Decause  $\delta(x_{i}) = h(x_{i}) - g(x_{i}) = y_{i} - y_{i} = 0.$   
So  $\langle h, h \rangle = \langle g, g \rangle + \langle \delta, \delta \rangle > \langle g, g \rangle$  Done!

## A proof of the smoothing spline solution

The proof is by contradiction and uses the interpolation result.

Let  $\hat{g}$  be the smoothing spline obtained as a linear combination of the kernel basis functions and possibly a linear or low order polynomial. This is found as a penalized smoother by plugging this form into the penalized least squares criterion and minimizing by ordinary calculus.

Suppose there is an h that has a smaller penalized least squares than  $\hat{g}$ .

Construct a g that interpolates h at the  $x_i$  and uses the same basis functions as  $\hat{g}$ . By the interpolation result we know that  $\langle h,h\rangle > \langle g,g\rangle$ and since both h and g have the same residuals sums of squares the penalized least squares criterion will be smaller for g. Thus , h can not be a minimizer.

Because  $\hat{g}$  however, has the same form as the interpolant and minimizes the criterion. Thus it must in fact be the minimizer over all functions in the Hilbert space.

## Choosing $\lambda$ by Cross-validation

Sequentially leave each observation out and predict it using the rest of the data. Find the  $\lambda$  that gives the best out of sample predictions.

Refitting the spline when each data point is omitted, and for a grid of  $\lambda$  values is computationally demanding.

Fortunately there is a shortcut.

The magic formula residual for  $g(x_i)$  having omitted  $y_i$ 

$$(y_i - \hat{g}_{-i}) = (y_i - \hat{g}_i)/(1 - A(\lambda))_{i,i}$$

This has a simple form because adding a data pair  $(x_i, \hat{g}_{-1})$  to the data does not change the estimate.

## **CV** and Generalized **CV** criterion

 $CV(\lambda)$ 

$$(1/n)\sum_{i=1}^{n}(y_i-\hat{g}_{-i})^2 = (1/n)\sum_{i=1}^{n}\frac{(y_i-\hat{g}_i)^2}{(1-A(\lambda))_{i,i})^2}$$

 $GCV(\lambda)$ 

$$(1/n)rac{\sum_{i=1}^n(y_i-\widehat{g}_i)^2}{(1-\operatorname{tr} A(\lambda)/n)^2}$$

Minimize CV or GCV over  $\lambda$  to determine a good value

## GCV for the ozone data

GCV( eff. degrees of freedom), the estimated surface



## GCV for the climate data

GCV( eff. degrees of freedom), the estimated surface



## Summary

We have formulated the curve/surface fitting problem as penalized least squares.

Splines treat estimating the entire curve but also have a finite basis related to a covariance function (reproducing kernel).

One can use CV or GCV to find the smoothing parameter.

## Thank you!

