Bayesian Interpretations of Regularization

Charlie Frogner

9.520 Class 10

March 12, 2012

The Plan

Regularized least squares maps $\{(x_i, y_i)\}_{i=1}^n$ to a function that minimizes the regularized loss:

$$f_{S} = \underset{f \in \mathcal{H}}{\operatorname{arg\,min}} \frac{1}{2} \sum_{i=1}^{n} (y_{i} - f(x_{i}))^{2} + \frac{\lambda}{2} ||f||_{\mathcal{H}}^{2}$$

Can we interpret RLS from a probabilistic point of view?

Some notation

- Training set: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}.$
- Inputs: $X = \{x_1, ..., x_n\}.$
- Labels: $Y = \{y_1, \dots, y_n\}.$
- Parameters: $\theta \in \mathbb{R}^p$.
- $p(Y|X, \theta)$ is the joint distribution over labels Y given inputs X and the parameters.

Where do probabilities show up?

$$\frac{1}{2} \sum_{i=1}^{n} V(y_i, f(x_i)) + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

becomes

$$p(\mathbf{Y}|f,\mathbf{X})\cdot p(f)$$

- Likelihood, a.k.a. noise model: p(Y|f, X).
 - Gaussian: $y_i \sim \mathcal{N}\left(f^*(x_i), \sigma_i^2\right)$
 - Poisson: $y_i \sim Pois(f^*(x_i))$
- Prior: p(f).

Where do probabilities show up?

$$\frac{1}{2} \sum_{i=1}^{n} V(y_i, f(x_i)) + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

becomes

$$p(\mathbf{Y}|f,\mathbf{X})\cdot p(f)$$

- Likelihood, a.k.a. noise model: p(Y|f, X).
 - Gaussian: $y_i \sim \mathcal{N}\left(f^*(x_i), \sigma_i^2\right)$
 - Poisson: $y_i \sim Pois(f^*(x_i))$
- **Prior**: *p*(*f*).



Estimation

The estimation problem:

- Given data $\{(x_i, y_i)\}_{i=1}^N$ and model $p(\mathbf{Y}|f, \mathbf{X}), p(f)$.
- Find a good f to explain data.

The Plan

- Maximum likelihood estimation for ERM
- MAP estimation for linear RLS
- MAP estimation for kernel RLS
- Transductive model
- Infinite dimensions get more complicated

Maximum likelihood estimation

- Given data $\{(x_i, y_i)\}_{i=1}^N$ and model $p(\mathbf{Y}|f, \mathbf{X}), p(f)$.
- A good f is one that maximizes p(Y|f, X).

Maximum likelihood and least squares

For least squares, noise model is:

$$y_i|f, \mathbf{x}_i \sim \mathcal{N}\left(f(\mathbf{x}_i), \sigma^2\right)$$

a.k.a.

$$\mathbf{Y}|f,\mathbf{X} \sim \mathcal{N}\left(f(\mathbf{X}),\sigma^2I\right)$$

So

$$p(\mathbf{Y}|f, \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\sum_{i=1}^{N} \frac{1}{\sigma^2} (y_i - f(x_i))^2\right\}$$

Maximum likelihood and least squares

For least squares, noise model is:

$$y_i|f, \mathbf{x}_i \sim \mathcal{N}\left(f(\mathbf{x}_i), \sigma^2\right)$$

a.k.a.

$$\mathbf{Y}|f,\mathbf{X}\sim\mathcal{N}\left(f(\mathbf{X}),\sigma^2I\right)$$

So

$$p(\mathbf{Y}|f, \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\sum_{i=1}^{N} \frac{1}{\sigma^2} (y_i - f(x_i))^2\right\}$$

Maximum likelihood and least squares

Maximum likelihood: maximize

$$p(\mathbf{Y}|f, \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\sum_{i=1}^{N} \frac{1}{\sigma^2} (y_i - f(x_i)))^2\right\}$$

Empirical risk minimization: minimize

$$\sum_{i=1}^{N}(y_i-f(x_i))^2$$

$$\sum_{i=1}^N (y_i - f(x_i))^2$$

$$e^{-\sum_{i=1}^{N} \frac{1}{\sigma^2} (y_i - f(x_i))^2}$$

RLS:

$$\arg\min_{f} \frac{1}{2} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

Is there a model of **Y** and f that yields RLS?

$$e^{-rac{1}{2\sigma_{arepsilon}^2}\left(\sum\limits_{i=1}^n(y_i-f(x_i))^2
ight)-rac{\lambda}{2}\|f\|_{\mathcal{H}}^2}$$

$$p(\mathbf{Y}|f,\mathbf{X}) \cdot p(f)$$

RLS:

$$\arg\min_{f} \frac{1}{2} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

Is there a model of **Y** and *f* that yields RLS?

$$e^{-rac{1}{2\sigma_{\varepsilon}^{2}}\left(\sum\limits_{i=1}^{n}(y_{i}-f(x_{i}))^{2}
ight)-rac{\lambda}{2}\|f\|_{\mathcal{H}}^{2}}$$

$$p(\mathbf{Y}|f,\mathbf{X}) \cdot p(f)$$



RLS:

$$\arg\min_{f} \frac{1}{2} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

Is there a model of **Y** and *f* that yields RLS?

$$e^{-\frac{1}{2\sigma_{\varepsilon}^2}\left(\sum\limits_{i=1}^n(y_i-f(x_i))^2\right)}\cdot e^{-\frac{\lambda}{2}\|f\|_{\mathcal{H}}^2}$$

$$p(\mathbf{Y}|f,\mathbf{X}) \cdot p(f)$$



RLS:

$$\arg\min_{f} \frac{1}{2} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \frac{\lambda}{2} ||f||_{\mathcal{H}}^2$$

Is there a model of **Y** and *f* that yields RLS?

$$e^{-\frac{1}{2\sigma_{\varepsilon}^2}\left(\sum\limits_{i=1}^n(y_i-f(x_i))^2\right)}\cdot e^{-\frac{\lambda}{2}\|f\|_{\mathcal{H}}^2}$$

$$p(\mathbf{Y}|f,\mathbf{X})\cdot p(f)$$

Posterior function estimates

- Given data $\{(x_i, y_i)\}_{i=1}^N$ and model $p(\mathbf{Y}|f, \mathbf{X}), p(f)$.
- Find a good f to explain data.

(If we can get $p(f|\mathbf{Y}, \mathbf{X})$)

Bayes least squares estimate:

$$\hat{f}_{BLS} = \mathbb{E}_{(f|\mathbf{X},\mathbf{Y})}[f]$$

i.e. the mean of the posterior.

MAP estimate:

$$\hat{f}_{MAP} = \arg\max_{f} p(f|\mathbf{X}, \mathbf{Y})$$

i.e. a mode of the posterior.



Posterior function estimates

- Given data $\{(x_i, y_i)\}_{i=1}^N$ and model $p(\mathbf{Y}|f, \mathbf{X}), p(f)$.
- Find a good *f* to explain data.

(If we can get $p(f|\mathbf{Y}, \mathbf{X})$)

Bayes least squares estimate:

$$\hat{f}_{BLS} = \mathbb{E}_{(f|\mathbf{X},\mathbf{Y})}[f]$$

i.e. the mean of the posterior.

MAP estimate:

$$\hat{f}_{MAP} = \arg\max_{f} p(f|\mathbf{X}, \mathbf{Y})$$

i.e. a mode of the posterior.



Posterior function estimates

- Given data $\{(x_i, y_i)\}_{i=1}^N$ and model $p(\mathbf{Y}|f, \mathbf{X}), p(f)$.
- Find a good *f* to explain data.

(If we can get $p(f|\mathbf{Y}, \mathbf{X})$)

Bayes least squares estimate:

$$\hat{f}_{\mathsf{BLS}} = \mathbb{E}_{(f|\mathbf{X},\mathbf{Y})}[f]$$

i.e. the mean of the posterior.

MAP estimate:

$$\hat{f}_{MAP} = \underset{f}{\operatorname{arg\,max}} \, p(f|\mathbf{X},\mathbf{Y})$$

i.e. a mode of the posterior.



How to find $p(f|\mathbf{Y}, \mathbf{X})$? Bayes' rule:

$$\rho(f|\mathbf{X}, \mathbf{Y}) = \frac{\rho(\mathbf{Y}|\mathbf{X}, f) \cdot \rho(f)}{\rho(\mathbf{Y}|\mathbf{X})}$$
$$= \frac{\rho(\mathbf{Y}|\mathbf{X}, f) \cdot \rho(f)}{\int \rho(\mathbf{Y}|\mathbf{X}, f) d\rho(f)}$$

When is this well-defined?

How to find $p(f|\mathbf{Y}, \mathbf{X})$? Bayes' rule:

$$\rho(f|\mathbf{X}, \mathbf{Y}) = \frac{\rho(\mathbf{Y}|\mathbf{X}, f) \cdot \rho(f)}{\rho(\mathbf{Y}|\mathbf{X})}$$
$$= \frac{\rho(\mathbf{Y}|\mathbf{X}, f) \cdot \rho(f)}{\int \rho(\mathbf{Y}|\mathbf{X}, f) d\rho(f)}$$

When is this well-defined?

Functions vs. parameters:

$$\mathcal{H} \cong \mathbb{R}^p$$

Represent functions in $\ensuremath{\mathcal{H}}$ by their coordinates w.r.t. a basis:

$$f \in \mathcal{H} \leftrightarrow \theta \in \mathbb{R}^p$$

Assume (for the moment): $ho < \infty$

Functions vs. parameters:

$$\mathcal{H}\cong\mathbb{R}^p$$

Represent functions in $\ensuremath{\mathcal{H}}$ by their coordinates w.r.t. a basis:

$$f \in \mathcal{H} \leftrightarrow \theta \in \mathbb{R}^p$$

Assume (for the moment): $p < \infty$

Mercer's theorem:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\mathbf{k}} \nu_{\mathbf{k}} \psi_{\mathbf{k}}(\mathbf{x}_i) \psi_{\mathbf{k}}(\mathbf{x}_j)$$

where $\nu_k \psi_k(\cdot) = \int K(\cdot, y) \psi_k(y) dy$ for all k. The functions $\{\sqrt{\nu_k} \psi_k(\cdot)\}$ form an *orthonormal basis* for \mathcal{H}_K . Let $\phi(\cdot) = [\sqrt{\nu_1} \psi_1(\cdot), \dots, \sqrt{\nu_p} \psi_p(\cdot)]$. Then:

$$\mathcal{H}_{K} = \{ \phi(\cdot)\theta | \theta \in \mathbb{R}^{p} \}$$

Prior on infinite-dimensional space

Problem: there's no such thing as

$$\theta \sim \mathcal{N}\left(\mathbf{0}, I\right)$$

when $\theta \in \mathbb{R}^{\infty}$!

Linear function:

$$f(x) = \langle x, \theta \rangle$$

Noise model:

$$\mathbf{Y}|\mathbf{X}, heta \sim \mathcal{N}\left(\mathbf{X} heta, \sigma_{arepsilon}^2 \mathbf{I}
ight)$$

Add a prior.

$$\theta \sim \mathcal{N}\left(\mathbf{0}, I\right)$$

Model:

$$\mathbf{Y}|\mathbf{X}, heta \sim \mathcal{N}\left(\mathbf{X} heta, \sigma_{arepsilon}^{2}\mathbf{I}
ight), \qquad heta \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}
ight)$$

Joint over **Y** and θ :

$$\left[\begin{array}{c} \mathbf{Y} \\ \boldsymbol{\theta} \end{array}\right] \sim \mathcal{N}\left(\left[\begin{array}{c} \mathbf{0} \\ \mathbf{0} \end{array}\right], \left[\begin{array}{cc} \mathbf{X}\mathbf{X}^T + \sigma_\varepsilon^2 \mathbf{I} & \mathbf{X} \\ \mathbf{X}^T & \mathbf{I} \end{array}\right]\right)$$

Condition on Y.



Posterior:

$$heta | \mathbf{X}, \mathbf{Y} \sim \mathcal{N}\left(\mu_{ heta | \mathbf{X}, \mathbf{Y}}, \Sigma_{ heta | \mathbf{X}, \mathbf{Y}}\right)$$

where

$$\begin{split} & \mu_{\theta|\mathbf{X},\mathbf{Y}} = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{Y} \\ & \boldsymbol{\Sigma}_{\theta|\mathbf{X},\mathbf{Y}} = I - \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{X} \end{split}$$

This is Gaussian, so

$$\hat{\theta}_{MAP}(\mathbf{Y}|\mathbf{X}) = \hat{\theta}_{BLS}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}^{T}(\mathbf{X}\mathbf{X}^{T} + \sigma_{\varepsilon}^{2}I)^{-1}\mathbf{Y}$$



Posterior:

$$heta | \mathbf{X}, \mathbf{Y} \sim \mathcal{N}\left(\mu_{ heta | \mathbf{X}, \mathbf{Y}}, \Sigma_{ heta | \mathbf{X}, \mathbf{Y}}\right)$$

where

$$\begin{split} & \mu_{\theta|\mathbf{X},\mathbf{Y}} = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{Y} \\ & \Sigma_{\theta|\mathbf{X},\mathbf{Y}} = I - \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{X} \end{split}$$

This is Gaussian, so

$$\hat{\theta}_{\textit{MAP}}(\mathbf{Y}|\mathbf{X}) = \hat{\theta}_{\textit{BLS}}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}^{T}(\mathbf{X}\mathbf{X}^{T} + \sigma_{\varepsilon}^{2}\mathbf{I})^{-1}\mathbf{Y}$$



Linear RLS as a MAP estimator

Model:

$$\mathbf{Y}|\mathbf{X}, heta \sim \mathcal{N}\left(\mathbf{X} heta, \sigma_{arepsilon}^{2} \mathbf{I}
ight), \qquad heta \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}
ight)$$

$$\hat{\theta}_{\textit{MAP}}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{Y}$$

Recall the linear RLS solution:

$$\hat{\theta}_{RLS}(\mathbf{Y}|\mathbf{X}) = \underset{\theta}{\arg\min} \frac{1}{2} \sum_{i=1}^{N} (y_i - \langle x_i, \theta \rangle)^2 + \frac{\lambda}{2} \|\theta\|^2$$
$$= \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \frac{\lambda}{2} I)^{-1} \mathbf{Y}$$

So what's λ ?



Linear RLS as a MAP estimator

Model:

$$\mathbf{Y}|\mathbf{X}, heta \sim \mathcal{N}\left(\mathbf{X} heta, \sigma_{arepsilon}^{2} \mathbf{I}
ight), \qquad heta \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}
ight)$$

$$\hat{\theta}_{\textit{MAP}}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{Y}$$

Recall the linear RLS solution:

$$\hat{\theta}_{RLS}(\mathbf{Y}|\mathbf{X}) = \underset{\theta}{\arg\min} \frac{1}{2} \sum_{i=1}^{N} (y_i - \langle x_i, \theta \rangle)^2 + \frac{\lambda}{2} \|\theta\|^2$$
$$= \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \frac{\lambda}{2} I)^{-1} \mathbf{Y}$$

So what's λ ?



Posterior for kernel RLS

Model for *linear* RLS:

$$\mathbf{Y}|\mathbf{X}, heta \sim \mathcal{N}\left(\mathbf{X} heta, \sigma_{arepsilon}^{2}\mathbf{I}
ight), \qquad heta \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}
ight)$$

Model for kernel RLS?

$$\mathbf{Y}|\mathbf{X},\theta \sim \mathcal{N}\left(\phi(\mathbf{X})\theta,\sigma_{\varepsilon}^{2}I\right), \qquad \theta \sim \mathcal{N}\left(\mathbf{0},I\right)$$

Then:

$$\hat{\theta}_{MAP}(\mathbf{Y}|\mathbf{X}) = \phi(\mathbf{X})^{T} (\phi(\mathbf{X})\phi(\mathbf{X})^{T} + \sigma_{\varepsilon}^{2} I)^{-1} \mathbf{Y}$$

Posterior for kernel RLS

Model for *linear* RLS:

$$\mathbf{Y}|\mathbf{X}, heta\sim\mathcal{N}\left(\mathbf{X} heta,\sigma_{arepsilon}^{2}\emph{\emph{I}}
ight),\qquad heta\sim\mathcal{N}\left(\mathbf{0},\emph{\emph{I}}
ight)$$

Model for kernel RLS?

$$\mathbf{Y}|\mathbf{X}, \theta \sim \mathcal{N}\left(\phi(\mathbf{X})\theta, \sigma_{\varepsilon}^{2}I\right), \qquad \theta \sim \mathcal{N}\left(0, I\right)$$

Then:

$$\hat{\theta}_{MAP}(\mathbf{Y}|\mathbf{X}) = \phi(\mathbf{X})^T (\phi(\mathbf{X})\phi(\mathbf{X})^T + \sigma_{\varepsilon}^2 I)^{-1} \mathbf{Y}$$

Posterior for kernel RLS

Model for *linear* RLS:

$$\mathbf{Y}|\mathbf{X}, heta\sim\mathcal{N}\left(\mathbf{X} heta,\sigma_{arepsilon}^{2}\emph{\emph{I}}
ight),\qquad heta\sim\mathcal{N}\left(0,\emph{\emph{I}}
ight)$$

Model for kernel RLS?

$$\mathbf{Y}|\mathbf{X}, \theta \sim \mathcal{N}\left(\phi(\mathbf{X})\theta, \sigma_{\varepsilon}^{2}I\right), \qquad \theta \sim \mathcal{N}\left(0, I\right)$$

Then:

$$\hat{ heta}_{MAP}(\mathbf{Y}|\mathbf{X}) = \phi(\mathbf{X})^T (\mathbf{K} + \sigma_{arepsilon}^2 \mathbf{I})^{-1} \mathbf{Y}$$

A quick recap

• Empirical risk minimization is ML.

$$ho(\mathbf{Y}|f,\mathbf{X}) \propto e^{-rac{1}{2}\sum_{i=1}^N(y_i-f(x_i))^2}$$

Linear RLS is MAP.

$$\rho(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - \langle x_i, \theta \rangle)^2} \cdot e^{-\frac{\lambda}{2} \theta^T \theta}$$

Kernel RLS is also MAP.

$$ho(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - f(x_i)^2 \cdot e^{-\frac{\lambda}{2} ||f||_{\mathcal{H}}^2}$$



A quick recap

Empirical risk minimization is ML.

$$ho(\mathbf{Y}|f,\mathbf{X}) \propto e^{-rac{1}{2}\sum_{i=1}^N(y_i-f(x_i))^2}$$

Linear RLS is MAP.

$$\rho(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - \langle \mathbf{x}_i, \theta \rangle)^2} \cdot e^{-\frac{\lambda}{2} \theta^T \theta}$$

Kernel RLS is also MAP.

$$ho(\mathbf{Y}, f | \mathbf{X}) \propto e^{-\frac{1}{2} \sum_{i=1}^{N} (y_i - f(x_i)^2 \cdot e^{-\frac{\lambda}{2}} \|f\|_{\mathcal{H}}^2}$$



A quick recap

Empirical risk minimization is ML.

$$ho(\mathbf{Y}|f,\mathbf{X}) \propto e^{-rac{1}{2}\sum_{i=1}^N(y_i-f(x_i))^2}$$

Linear RLS is MAP.

$$\rho(\mathbf{Y},f|\mathbf{X}) \propto e^{-\frac{1}{2}\sum_{i=1}^{N}(y_i - \langle x_i,\theta \rangle)^2} \cdot e^{-\frac{\lambda}{2}\theta^T\theta}$$

Kernel RLS is also MAP.

$$ho(\mathbf{Y},f|\mathbf{X}) \propto e^{-rac{1}{2}\sum_{i=1}^{N}(y_i-f(x_i)^2}\cdot e^{-rac{\lambda}{2}\|f\|_{\mathcal{H}}^2}$$



Idea: Forget about estimating θ (i.e. f).

Instead: Estimate predicted outputs

$$\mathbf{Y}^* = [y_1^*, \dots, y_M^*]^T$$

at test inputs

$$\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_M^*]^T$$

Need the joint distribution over \mathbf{Y}^* and \mathbf{Y} .



Say **Y*** and **Y** are *jointly Gaussian*:

$$\left[\begin{array}{c} \boldsymbol{Y} \\ \boldsymbol{Y}^* \end{array}\right] = \mathcal{N}\left(\left[\begin{array}{c} \boldsymbol{0} \\ \boldsymbol{0} \end{array}\right], \left[\begin{array}{cc} \boldsymbol{\Lambda}_{\boldsymbol{Y}} & \boldsymbol{\Lambda}_{\boldsymbol{Y}\boldsymbol{Y}^*} \\ \boldsymbol{\Lambda}_{\boldsymbol{Y}^*\boldsymbol{Y}} & \boldsymbol{\Lambda}_{\boldsymbol{Y}^*} \end{array}\right]\right)$$

Want: kernel RLS.

General form for the posterior:

$$\mathbf{Y}^*|\mathbf{X},\mathbf{Y}\sim\mathcal{N}\left(\boldsymbol{\mu}_{\mathbf{Y}*|\mathbf{X},\mathbf{Y}},\boldsymbol{\Sigma}_{\mathbf{Y}^*|\mathbf{X},\mathbf{Y}}\right)$$

where

$$\begin{split} &\mu_{\mathsf{Y}^*|\mathbf{X},\mathbf{Y}} = \Lambda_{\mathsf{Y}\mathsf{Y}^*}^T \Lambda_{\mathsf{Y}}^{-1} \mathbf{Y} \\ &\Sigma_{\mathsf{Y}^*|\mathbf{X},\mathbf{Y}} = \Lambda_{\mathsf{Y}^*} - \Lambda_{\mathsf{Y}\mathsf{Y}^*}^T \Lambda_{\mathsf{Y}}^{-1} \Lambda_{\mathsf{Y}\mathsf{Y}^*} \end{split}$$



Say Y* and Y are jointly Gaussian:

$$\left[\begin{array}{c} \boldsymbol{Y} \\ \boldsymbol{Y}^* \end{array}\right] = \mathcal{N}\left(\left[\begin{array}{c} \boldsymbol{0} \\ \boldsymbol{0} \end{array}\right], \left[\begin{array}{cc} \boldsymbol{\Lambda}_{\boldsymbol{Y}} & \boldsymbol{\Lambda}_{\boldsymbol{Y}\boldsymbol{Y}^*} \\ \boldsymbol{\Lambda}_{\boldsymbol{Y}^*\boldsymbol{Y}} & \boldsymbol{\Lambda}_{\boldsymbol{Y}^*} \end{array}\right]\right)$$

Want: kernel RLS.

General form for the posterior:

$$\mathbf{Y}^*|\mathbf{X},\mathbf{Y}\sim\mathcal{N}\left(\boldsymbol{\mu}_{\mathbf{Y}*|\mathbf{X},\mathbf{Y}},\boldsymbol{\Sigma}_{\mathbf{Y}^*|\mathbf{X},\mathbf{Y}}\right)$$

where

$$\mu_{\mathbf{Y}^*|\mathbf{X},\mathbf{Y}} = \Lambda_{\mathbf{Y}^*}^T \Lambda_{\mathbf{Y}}^{-1} \mathbf{Y}$$

$$\Sigma_{\mathbf{Y}^*|\mathbf{X},\mathbf{Y}} = \Lambda_{\mathbf{Y}^*} - \Lambda_{\mathbf{Y}^*}^T \Lambda_{\mathbf{Y}^*}^{-1} \Lambda_{\mathbf{Y}^*}$$

Set
$$\Lambda_{\mathbf{Y}} = K(\mathbf{X}, \mathbf{X}) + \sigma^2 I$$
, $\Lambda_{\mathbf{Y}Y^*} = K(\mathbf{X}, X^*)$, $\Lambda_{Y^*} = K(X^*, X^*)$.

Posterior:

$$\mathbf{Y}^* | \mathbf{X}, \mathbf{Y} \sim \mathcal{N}\left(\mu_{\mathbf{Y}_* | \mathbf{X}, \mathbf{Y}}, \Sigma_{\mathbf{Y}^* | \mathbf{X}, \mathbf{Y}}\right)$$

where

$$\begin{split} &\mu_{\mathsf{Y}^*|\mathbf{X},\mathbf{Y}} = K(X^*,\mathbf{X})(K(\mathbf{X},\mathbf{X}) + \sigma^2 I)^{-1}\mathbf{Y} \\ &\Sigma_{\mathsf{Y}^*|\mathbf{X},\mathbf{Y}} = K(X^*,X^*) - K(X^*,\mathbf{X})(K(\mathbf{X},\mathbf{X}) + \sigma^2 I)^{-1}K(X,X^*) \end{split}$$

So:
$$\hat{Y}_{MAP}^* = \hat{f}_{RLS}(X^*)$$
.



Model:

$$\left[\begin{array}{c} \mathbf{Y} \\ \mathbf{Y}^* \end{array}\right] = \mathcal{N}\left(\left[\begin{array}{c} \mathbf{0} \\ \mathbf{0} \end{array}\right], \left[\begin{array}{cc} \mathbf{K}(\mathbf{X},\mathbf{X}) + \sigma_\varepsilon^2 \mathbf{I} & \mathbf{K}(\mathbf{X},\mathbf{X}^*) \\ \mathbf{K}(\mathbf{X}^*,\mathbf{X}) & \mathbf{K}(\mathbf{X}^*,\mathbf{X}^*) \end{array}\right]\right)$$

MAP estimate (posterior mean) = RLS function at every point x^* , regardless of dim \mathcal{H}_K .

Are the prior and posterior (*on points*!) consistent with a distribution on \mathcal{H}_K ?

Strictly speaking, θ and f don't come into play here at all:

Have: $p(Y^*|X,Y)$ Do not have: $p(\theta|X,Y)$ or p(f|X,Y)

But, if \mathcal{H}_K is finite dimensional, the joint over Y and Y* is consistent with:

- $\mathbf{Y} = f(\mathbf{X}) + \varepsilon$,
- $Y^* = f(X)$, and
- f ∈ H_K is a random trajectory from a Gaussian process over the domain, with mean μ and covariance K.
- (Ergo, people call this "Gaussian process regression.")
 (Also "Kriging," because of a guy.)



Strictly speaking, θ and f don't come into play here at all:

Have:
$$p(Y^*|\mathbf{X}, \mathbf{Y})$$

Do not have: $p(\theta|\mathbf{X}, \mathbf{Y})$ or $p(f|\mathbf{X}, \mathbf{Y})$

But, if \mathcal{H}_K is finite dimensional, the joint over Y and Y* is consistent with:

- $\mathbf{Y} = f(\mathbf{X}) + \varepsilon$,
- $Y^* = f(X)$, and
- $f \in \mathcal{H}_K$ is a random trajectory from a **Gaussian process** over the domain, with mean μ and covariance K.
- (Ergo, people call this "Gaussian process regression.")
 (Also "Kriging," because of a guy.)



Recap redux

 Empirical risk minimization is the maximum likelihood estimator when:

$$\mathbf{y} = \mathbf{x}^T \theta + \varepsilon$$

Linear RLS is the MAP estimator when:

$$\mathbf{y} = \mathbf{x}^T \theta + \varepsilon, \qquad \theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

• Kernel RLS is the MAP estimator when:

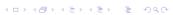
$$\mathbf{y} = \phi(\mathbf{x})^T \theta + \varepsilon, \qquad \theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

in finite dimensional \mathcal{H}_K .

• Kernel RLS is the MAP estimator at points when:

$$\left[\begin{array}{c} \mathbf{Y} \\ \mathbf{Y}^* \end{array}\right] = \mathcal{N}\left(\left[\begin{array}{c} \mu_{\mathbf{Y}} \\ \mu_{\mathbf{Y}^*} \end{array}\right], \left[\begin{array}{cc} \mathbf{K}(\mathbf{X},\mathbf{X}) + \sigma_{\varepsilon}^2 \mathbf{I} & \mathbf{K}(\mathbf{X},\mathbf{X}^*) \\ \mathbf{K}(\mathbf{X}^*,\mathbf{X}) & \mathbf{K}(\mathbf{X}^*,\mathbf{X}^*) \end{array}\right]\right)$$

in possibly infinite dimensional \mathcal{H}_K .



Is this useful in practice?

- Want confidence intervals + believe the posteriors are meaningful = yes
- Maybe other reasons?

