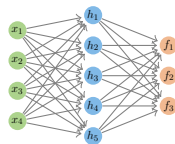# UNIVERSITY OF CAMBRIDGE

# Neural Nets, $\mathcal{GP}$s, and where the kernel lives

Paul Rubenstein[1][2]     Matthias Bauer[1][2]
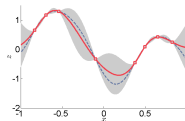
[1] University of Cambridge

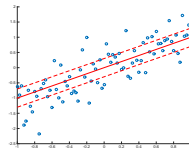[2] Max-Planck Institute for Intelligent Systems, Tübingen

12th November 2015

**Priors over infinite NN = $\mathcal{GP}$**



**Relationship between Kernel Ridge Regression and GPs**



**Support Vector Regression and GPs**



UNIVERSITY OF
CAMBRIDGE

# Acknowledgements/References

This talk is based mostly on the following:

- Arthur Gretton's course on RKHS theory: `http://www.gatsby.ucl.ac.uk/~gretton/coursefiles/rkhscourse.html`
- Bishop's Pattern Recognition and Machine Learning
- Stulp and Sigaud, Many regression algorithms, one unified model: A review

# Ordinary Least Squares (OLS) Linear Regression

**Problem:**

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \dots, N\}$ with $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^p$ and $y_i \in \mathcal{Y} = \mathbb{R}$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ that explains* $\mathcal{D}$.

# Ordinary Least Squares (OLS) Linear Regression

**Problem:**

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^p$ and $y_i \in \mathcal{Y} = \mathbb{R}$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ that explains* $\mathcal{D}$.

**An approach:**

- Assume $f$ is linear: $f(\mathbf{x}) = \mathbf{x}^\mathsf{T}\beta$ for some $\beta$
- Choose $\beta$ to minimise the sum of squared errors.

# Ordinary Least Squares (OLS) Linear Regression

**Problem:**

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^p$ and $y_i \in \mathcal{Y} = \mathbb{R}$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ that explains* $\mathcal{D}$.

**An approach:**

- Assume $f$ is linear: $f(\mathbf{x}) = \mathbf{x}^\intercal \beta$ for some $\beta$
- Choose $\beta$ to minimise the sum of squared errors.

Writing $X = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N)^\intercal$ and $Y = (y_1, \ldots, y_N)^\intercal$, we wish to minimise

$$L(\beta) = (Y - X\beta)^\intercal (Y - X\beta)$$

$$L(\beta) = (Y - X\beta)^{\mathsf{T}}(Y - X\beta)$$
$$= Y^{\mathsf{T}}Y - 2\beta^{\mathsf{T}}X^{\mathsf{T}}Y + \beta^{\mathsf{T}}X^{\mathsf{T}}X\beta$$
$$\implies \frac{dL}{d\beta} = -2X^{\mathsf{T}}Y + 2X^{\mathsf{T}}X\beta$$

So $\frac{dL}{d\beta} = 0 \implies \beta = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}Y$ if $(X^{\mathsf{T}}X)^{-1}$ exists.

# OLS (cont)

Two problems.

1. What if $X^\mathsf{T} X$ is not invertible?
2. What if $y$ is not well approximated by a linear function of $\mathbf{x}$?

# OLS (cont)

Two problems.

1. What if $X^\mathsf{T} X$ is not invertible?
2. What if $y$ is not well approximated by a linear function of $\mathbf{x}$?

Solutions:

1. Eigenvalues of $X^\mathsf{T} X$ are always $\geq 0$
   $\implies X^\mathsf{T} X + \lambda I$ invertible for $\lambda > 0$... why?
2. Can replace $\mathbf{x}$ with $\phi(\mathbf{x})$, where $\phi : \mathcal{X} \longrightarrow \mathbb{R}^p$.
   Write $\Phi := \phi(X)$

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathbb{R}^p$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ so that $f \circ \phi$ explains* $\mathcal{D}$.

# Ridge Regression in Feature Space

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathbb{R}^p$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ so that $f \circ \phi$ explains* $\mathcal{D}$.

**An approach**:

- Assume $f$ is linear: $f(\phi(\mathbf{x})) = \phi(\mathbf{x})^\mathsf{T} \beta$ for some $\beta$
- If $p$ is large compared to $N$ then we may overfit
- So choose $\beta$ to minimise the sum of squared errors plus complexity penalty.

$$L(\beta) = \sum_i (f(\phi(\mathbf{x}_i)) - y_i)^2 + \lambda \|\beta\|^2$$

$$= (Y - \Phi\beta)^\mathsf{T}(Y - \Phi\beta) + \lambda\beta^\mathsf{T}\beta$$

# Ridge Regression in Feature Space

$$L(\beta) = (Y - \Phi\beta)^{\mathsf{T}}(Y - \Phi\beta) + \lambda\beta^{\mathsf{T}}\beta$$

$$\frac{dL}{d\beta} = 0 \implies 0 = -2\Phi^{\mathsf{T}}Y + 2\Phi^{\mathsf{T}}\Phi\beta + 2\lambda\beta$$

$$\implies \beta = (\Phi^{\mathsf{T}}\Phi + \lambda_p I)^{-1}\Phi^{\mathsf{T}}Y$$

$$\beta = (\Phi^\intercal \Phi + \lambda_p I)^{-1} \Phi^\intercal Y$$

$$\beta = (\Phi^\intercal \Phi + \lambda_p I)^{-1} \Phi^\intercal Y$$

Two observations:

- $\Phi^\intercal \Phi$ is NOT the Gram matrix

$$\beta = (\Phi^{\mathsf{T}}\Phi + \lambda_p I)^{-1}\Phi^{\mathsf{T}}Y$$

Two observations:

- $\Phi^{\mathsf{T}}\Phi$ is NOT the Gram matrix
- $(\Phi^{\mathsf{T}}\Phi + \lambda I_p)\Phi^{\mathsf{T}} = \Phi^{\mathsf{T}}\Phi\Phi^{\mathsf{T}} + \lambda\Phi^{\mathsf{T}} = \Phi^{\mathsf{T}}(\Phi\Phi^{\mathsf{T}} + \lambda I_N)$

# Ridge Regression in Feature Space

$$\beta = (\Phi^\mathsf{T}\Phi + \lambda_p I)^{-1}\Phi^\mathsf{T}Y$$

Two observations:

- $\Phi^\mathsf{T}\Phi$ is NOT the Gram matrix
- $(\Phi^\mathsf{T}\Phi + \lambda I_p)\Phi^\mathsf{T} = \Phi^\mathsf{T}\Phi\Phi^\mathsf{T} + \lambda\Phi^\mathsf{T} = \Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)$

All eigenvalues of $\Phi^\mathsf{T}\Phi$ and $\Phi\Phi^\mathsf{T}$ are $\geq 0$ and so both bracketed expressions are invertible. Thus

$$\Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1} = (\Phi^\mathsf{T}\Phi + \lambda I_p)^{-1}\Phi^\mathsf{T}$$

# Regularised feature-mapped regression

So instead we can write

$$\beta = \Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1} Y$$
$$\implies f(\mathbf{x}_*) = \phi(\mathbf{x}_*)^\mathsf{T}\Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1} Y$$

# Regularised feature-mapped regression

So instead we can write

$$\beta = \Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1}Y$$
$$\implies f(\mathbf{x}_*) = \phi(\mathbf{x}_*)^\mathsf{T}\Phi^\mathsf{T}(\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1}Y$$

**Some reasons this might be good**:

- If $p > N$ then the matrix inversion takes $O(N^3)$ operations compared to $O(p^3)$
- $\phi(\mathbf{x})$ only ever appears as an inner product - so might not need to explicitly represent $\phi$

## Definition (Kernel)

A function $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ is a *kernel* if it is symmetric and if, for any $x_1, \ldots, x_n \in \mathcal{X}$, the matrix $K$ with entries $K_{ij} = k(x_i, x_j)$ is positive semi-definite.

# Brief Introduction to RKHS Theory

## Definition (Kernel)

A function $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ is a *kernel* if it is symmetric and if, for any $x_1, \ldots, x_n \in \mathcal{X}$, the matrix $K$ with entries $K_{ij} = k(x_i, x_j)$ is positive semi-definite.

$K$ positive semi-definite $\iff a^\mathsf{T} K a \geq 0$ for any $a \in \mathbb{R}^n$.

## Definition (Kernel)

A function $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ is a *kernel* if it is symmetric and if, for any $x_1, \ldots, x_n \in \mathcal{X}$, the matrix $K$ with entries $K_{ij} = k(x_i, x_j)$ is positive semi-definite.

$K$ positive semi-definite $\iff$ $a^\mathsf{T} K a \geq 0$ for any $a \in \mathbb{R}^n$.

**Example**: Let $\phi : \mathcal{X} \longrightarrow \mathcal{H}$ be any map into a Hilbert space, then $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$ is a kernel.

- **Symmetry**: inherited from $\langle ., . \rangle$
- **+ve semidefinite**:
  $a^\mathsf{T} K a = \sum_{ij} \langle a_i \phi(x_i), a_j \phi(x_j) \rangle_{\mathcal{H}} = \| \sum_i a_i \phi(x_i) \|_{\mathcal{H}}^2 \geq 0$

**UNIVERSITY OF CAMBRIDGE**

## Definition (Reproducing Kernel Hilbert Space)

Let $\mathcal{H}$ be a Hilbert space of functions $\mathcal{X} \longrightarrow \mathbb{R}$. We say that $\mathcal{H}$ is a Reproducing Kernel Hilbert Space (RKHS) if the evaluation operators $\delta_x : \mathcal{H} \longrightarrow \mathbb{R}, f \mapsto f(x)$ are continuous for all $x \in \mathcal{X}$

# Brief Introduction to RKHS Theory

## Definition (Reproducing Kernel Hilbert Space)

Let $\mathcal{H}$ be a Hilbert space of functions $\mathcal{X} \longrightarrow \mathbb{R}$. We say that $\mathcal{H}$ is a Reproducing Kernel Hilbert Space (RKHS) if the evaluation operators $\delta_x : \mathcal{H} \longrightarrow \mathbb{R}, f \mapsto f(x)$ are continuous for all $x \in \mathcal{X}$

$\delta_x$ continuous means...

- convergence in norm of a sequence of functions implies pointwise convergence at every point so functions are 'smooth'

# Brief Introduction to RKHS Theory

## Definition (Reproducing Kernel Hilbert Space)

Let $\mathcal{H}$ be a Hilbert space of functions $\mathcal{X} \longrightarrow \mathbb{R}$. We say that $\mathcal{H}$ is a Reproducing Kernel Hilbert Space (RKHS) if the evaluation operators $\delta_x : \mathcal{H} \longrightarrow \mathbb{R}, f \mapsto f(x)$ are continuous for all $x \in \mathcal{X}$

$\delta_x$ continuous means...

- convergence in norm of a sequence of functions implies pointwise convergence at every point so functions are 'smooth'
- by Riesz, there exists a unique $\phi_x \in \mathcal{H}$ such that $f(x) = \langle f, \phi_x \rangle$ for all $f \in \mathcal{H}$.

# Brief Introduction to RKHS Theory

## Definition (Reproducing Kernel Hilbert Space)

Let $\mathcal{H}$ be a Hilbert space of functions $\mathcal{X} \longrightarrow \mathbb{R}$. We say that $\mathcal{H}$ is a Reproducing Kernel Hilbert Space (RKHS) if the evaluation operators $\delta_x : \mathcal{H} \longrightarrow \mathbb{R}, f \mapsto f(x)$ are continuous for all $x \in \mathcal{X}$

$\delta_x$ continuous means...

- convergence in norm of a sequence of functions implies pointwise convergence at every point so functions are 'smooth'
- by Riesz, there exists a unique $\phi_x \in \mathcal{H}$ such that $f(x) = \langle f, \phi_x \rangle$ for all $f \in \mathcal{H}$.

We call $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}, (x, x') \mapsto \langle \phi_x, \phi_{x'} \rangle$ the (unique) Reproducing Kernel of $\mathcal{H}$

Summary:

- An RKHS on a base set $\mathcal{X}$ is just[1] a set of functions $\mathcal{X} \longrightarrow \mathbb{R}$
- Given an RKHS, we can construct a kernel on $\mathcal{X}$

---

[1]with some previously mentioned caveats

Summary:

- An RKHS on a base set $\mathcal{X}$ is just[1] a set of functions $\mathcal{X} \longrightarrow \mathbb{R}$
- Given an RKHS, we can construct a kernel on $\mathcal{X}$

Remarkably, the converse holds.

### Theorem (Moore-Aronszajn)

*Suppose that $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ is a kernel. Then there exists an RKHS $\mathcal{H}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathcal{H}$ such that*
*$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$*

---

[1]with some previously mentioned caveats

UNIVERSITY OF
CAMBRIDGE

# Brief Introduction to RKHS Theory

Summary:

- An RKHS on a base set $\mathcal{X}$ is just[1] a set of functions $\mathcal{X} \longrightarrow \mathbb{R}$
- Given an RKHS, we can construct a kernel on $\mathcal{X}$

Remarkably, the converse holds.

## Theorem (Moore-Aronszajn)

*Suppose that $k : \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ is a kernel. Then there exists an RKHS $\mathcal{H}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathcal{H}$ such that*
$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$

- $\mathcal{H}$ is the smallest Hilbert space containing each $k(\cdot, x)$
- properties of functions determined through properties of $k(\cdot, x)$

---

[1]with some previously mentioned caveats

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$
- Want to infer a function $f : \mathcal{X} \longrightarrow \mathbb{R}$ so that $f$ explains* $\mathcal{D}$.

# Kernel Ridge Regression

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \dots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$
- Want to infer a function $f : \mathcal{X} \longrightarrow \mathbb{R}$ so that $f$ explains* $\mathcal{D}$.

**An approach**:

- Pick a kernel $k$ such that the functions $k(\cdot, x)$ are 'good'
- Consider the RKHS $\mathcal{H}$ corresponding to $k$
- Find the $f \in \mathcal{H}$ that minimises empirical squared error (with penalty for complexity)

$$\underset{f \in \mathcal{H}}{\arg\min} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\underset{f \in \mathcal{H}}{\arg\min} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How do we find the argmin?

# Kernel Ridge Regression

$$\arg\min_{f \in \mathcal{H}} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How do we find the argmin? Answer:

## Theorem (Representer theorem)

*The solution $f_*$ to the above problem lies in the subspace of $\mathcal{H}$ spanned by the set $\{k(\cdot, x_i) | i = 1, \ldots, N\}$. ie*
*$f_* = \sum_i \alpha_i k(\cdot, x_i)$ for some coefficients $\alpha_i$*

UNIVERSITY OF
CAMBRIDGE

# Kernel Ridge Regression

$$\arg\min_{f \in \mathcal{H}} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How do we find the argmin? Answer:

## Theorem (Representer theorem)

*The solution $f_*$ to the above problem lies in the subspace of $\mathcal{H}$ spanned by the set $\{k(\cdot, x_i) | i = 1, \ldots, N\}$. ie*
*$f_* = \sum_i \alpha_i k(\cdot, x_i)$ for some coefficients $\alpha_i$*

$$\arg\min_{\alpha \in \mathbb{R}^N} \sum_i (f_\alpha(\mathbf{x}_i) - y_i)^2 + \lambda \|f_\alpha\|_{\mathcal{H}}^2$$

# Kernel Ridge Regression

Proof:

- Let $f \in \mathcal{H}$
- Let $f_s$ be the projection of $f$ onto $\mathrm{span}\{k(\cdot, x_i)\}$
- Let $f_\perp = f - f_s \perp \mathrm{span}\{k(\cdot, x_i)\}$

# Kernel Ridge Regression

<u>Proof</u>:

- Let $f \in \mathcal{H}$
- Let $f_s$ be the projection of $f$ onto $\mathrm{span}\{k(\cdot, x_i)\}$
- Let $f_\perp = f - f_s \perp \mathrm{span}\{k(\cdot, x_i)\}$

We show that $f_s$ is better than $f$ in the sense that:

- The loss function is the same: $(f_s(x) - y)^2 = (f(x) - y)^2$
- The complexity penalty is smaller: $\|f_s\|_{\mathcal{H}}^2 \leq \|f\|_{\mathcal{H}}^2$

# Kernel Ridge Regression

For each term in the loss function we have:

$$
\begin{aligned}
(f(x_i) - y_i)^2 &= (f_s(x_i) + f_\perp(x_i) - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i)\rangle + \langle f_\perp, k(\cdot, x_i)\rangle - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i)\rangle - y_i)^2 \\
&= (f_s(x_i) - y_i)^2
\end{aligned}
$$

UNIVERSITY OF
CAMBRIDGE

# Kernel Ridge Regression

For each term in the loss function we have:

$$
\begin{aligned}
(f(x_i) - y_i)^2 &= (f_s(x_i) + f_\perp(x_i) - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i)\rangle + \langle f_\perp, k(\cdot, x_i)\rangle - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i)\rangle - y_i)^2 \\
&= (f_s(x_i) - y_i)^2
\end{aligned}
$$

Considering the complexity penalty:

$$
\begin{aligned}
\|f\|_{\mathcal{H}}^2 &= \|f_s + f_\perp\|_{\mathcal{H}}^2 \\
&= \|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}^2 \geq \|f_s\|_{\mathcal{H}}^2
\end{aligned}
$$

For each term in the loss function we have:

$$
\begin{aligned}
(f(x_i) - y_i)^2 &= (f_s(x_i) + f_\perp(x_i) - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i) \rangle + \langle f_\perp, k(\cdot, x_i) \rangle - y_i)^2 \\
&= (\langle f_s, k(\cdot, x_i) \rangle - y_i)^2 \\
&= (f_s(x_i) - y_i)^2
\end{aligned}
$$

Considering the complexity penalty:

$$
\begin{aligned}
\|f\|_{\mathcal{H}}^2 &= \|f_s + f_\perp\|_{\mathcal{H}}^2 \\
&= \|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}^2 \geq \|f_s\|_{\mathcal{H}}^2
\end{aligned}
$$

So $f_s$ is better than $f$! Thus optimal $f_*$ must lie in $\mathrm{span}\{k(\cdot, x_i)\}$

# Kernel Ridge Regression

Writing $f = \sum_j \alpha_j k(\cdot, x_j)$, we wish to minimise the following quantity over $\alpha$:

$$
\begin{aligned}
L(\alpha) &= \sum_i (f(x_i) - y_i)^2 + \lambda \|f\|^2 \\
&= \sum_i (\sum_j \langle \alpha_j k(\cdot, x_j), k(\cdot, x_i) \rangle - y_i)^2 + \lambda \langle f, f \rangle \\
&= \sum_i ((K\alpha)_i - y_i)^2 + \lambda \sum_{ij} \langle \alpha_i k(\cdot, x_i), \alpha_j k(\cdot, x_j) \rangle \\
&= (K\alpha - Y)^\mathsf{T} (K\alpha - Y) + \lambda \alpha^\mathsf{T} K \alpha
\end{aligned}
$$

# Kernel Ridge Regression

Differentiating with respect to $\alpha$ yields

$$\frac{dL}{d\alpha} = 2KK\alpha - 2KY + 2\lambda K\alpha$$
$$= 2K(K\alpha - Y + \lambda\alpha)$$
$$= 2K((K + \lambda I_N)\alpha - Y) = 0$$
$$\implies \alpha = (K + \lambda I_N)^{-1}Y$$

# Kernel Ridge Regression

Differentiating with respect to $\alpha$ yields

$$\begin{aligned}
\frac{dL}{d\alpha} &= 2KK\alpha - 2KY + 2\lambda K\alpha \\
&= 2K(K\alpha - Y + \lambda\alpha) \\
&= 2K((K + \lambda I_N)\alpha - Y) = 0 \\
\implies \alpha &= (K + \lambda I_N)^{-1}Y
\end{aligned}$$

For a new point $x_*$, writing $\mathbf{k}$ to be the vector with $\mathbf{k}_i = k(x_*, x_i)$ we see that
$f(x_*) = \sum_i \alpha_i k(x_*, x_i) = \mathbf{k}^\mathsf{T}(K + \lambda I_N)^{-1}Y$

# Solution is the same as before

Old solution:

$$f(\mathbf{x}_*) = \phi(\mathbf{x}_*)^\mathsf{T} \Phi^\mathsf{T} (\Phi\Phi^\mathsf{T} + \lambda I_N)^{-1} Y$$

New solution:

$$f(x_*) = \mathbf{k}^\mathsf{T} (K + \lambda I_N)^{-1} Y$$

▶ If we look 'inside' the $\mathbf{k}$ and $K$, we see that these are the same.

- Starting with linear regression, we have derived Kernel Ridge Regression.

# Summary so far

▶ Starting with linear regression, we have derived Kernel Ridge Regression.

▶ Crucial idea 1: Regulariser $\lambda$ let us write all computations in terms of inner products between feature mapped observations.

# Summary so far

- ▶ Starting with linear regression, we have derived Kernel Ridge Regression.

- ▶ Crucial idea 1: Regulariser $\lambda$ let us write all computations in terms of inner products between feature mapped observations.

- ▶ Crucial idea 2: Representer theorem $\implies$ can project infinite dimensional optimisation problem to finite dimensional space

# Summary so far

- Starting with linear regression, we have derived Kernel Ridge Regression.
- Crucial idea 1: Regulariser $\lambda$ let us write all computations in terms of inner products between feature mapped observations.
- Crucial idea 2: Representer theorem $\implies$ can project infinite dimensional optimisation problem to finite dimensional space

Diferent approach?

- Motivation to use regulariser was to prevent overfitting
- Could instead adopt a Bayesian approach

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathbb{R}^p$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ so that $f \circ \phi$ explains $\mathcal{D}$.

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \mathbb{R}$ and feature map $\phi : \mathcal{X} \longrightarrow \mathbb{R}^p$
- Want to infer a function $f : \mathbb{R}^p \longrightarrow \mathbb{R}$ so that $f \circ \phi$ explains $\mathcal{D}$.

**An approach**:

- Assume $f$ is linear: $f(\phi(\mathbf{x})) = \phi(\mathbf{x})^\mathsf{T}\alpha$ for some $\alpha$
- Place prior over $\alpha$, add noise and perform Bayesian inference

$$y(\mathbf{x}) = \phi(\mathbf{x})^\mathsf{T}\mathbf{w} + \epsilon \qquad \mathbf{w} \sim \mathcal{N}(\mathbf{w}|0, \sigma_w^2 I), \quad \epsilon \sim \mathcal{N}(\epsilon|0, \sigma_\epsilon^2)$$

# GP Regression

$\mathbf{w}, \epsilon$ Gaussian $\implies y$ is Gaussian with

$$\mathbb{E}(y(\mathbf{x})) = \phi(\mathbf{x})^{\mathsf{T}}\mathbb{E}(\mathbf{w}) + \mathbb{E}(\epsilon) = 0$$
$$\mathrm{Cov}(y(\mathbf{x}), y(\mathbf{x}')) = \sigma_w^2 \phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}') + \delta_{\mathbf{x}=\mathbf{x}'}\sigma_\epsilon^2$$

UNIVERSITY OF
CAMBRIDGE

$\mathbf{w}, \epsilon$ Gaussian $\implies y$ is Gaussian with

$$\mathbb{E}(y(\mathbf{x})) = \phi(\mathbf{x})^\mathsf{T}\mathbb{E}(\mathbf{w}) + \mathbb{E}(\epsilon) = 0$$
$$\mathrm{Cov}(y(\mathbf{x}), y(\mathbf{x}')) = \sigma_w^2 \phi(\mathbf{x})^\mathsf{T}\phi(\mathbf{x}') + \delta_{\mathbf{x}=\mathbf{x}'}\sigma_\epsilon^2$$

- Write $K$ for the matrix with $K_{ij} = \phi(\mathbf{x}_i)^\mathsf{T}\phi(\mathbf{x}_j)$
- $\mathbf{k}$ for the vector with $\mathbf{k}_i = \phi(\mathbf{x}_*)^\mathsf{T}\phi(\mathbf{x}_i)$
- $c = \phi(\mathbf{x}_*)^\mathsf{T}\phi(\mathbf{x}_*)$
- $\mathbf{y} = (y_1, \ldots, y_N, y_*)^\mathsf{T}$

$\mathbf{w}, \epsilon$ Gaussian $\implies y$ is Gaussian with

$$\mathbb{E}(y(\mathbf{x})) = \phi(\mathbf{x})^{\mathsf{T}}\mathbb{E}(\mathbf{w}) + \mathbb{E}(\epsilon) = 0$$
$$\mathrm{Cov}(y(\mathbf{x}), y(\mathbf{x}')) = \sigma_w^2\phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}') + \delta_{\mathbf{x}=\mathbf{x}'}\sigma_\epsilon^2$$

- ▸ Write $K$ for the matrix with $K_{ij} = \phi(\mathbf{x}_i)^{\mathsf{T}}\phi(\mathbf{x}_j)$
- ▸ $\mathbf{k}$ for the vector with $\mathbf{k}_i = \phi(\mathbf{x}_*)^{\mathsf{T}}\phi(\mathbf{x}_i)$
- ▸ $c = \phi(\mathbf{x}_*)^{\mathsf{T}}\phi(\mathbf{x}_*)$
- ▸ $\mathbf{y} = (y_1, \ldots, y_N, y_*)^{\mathsf{T}}$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{y}|0, \sigma_w^2 \begin{pmatrix} K + \frac{\sigma_\epsilon^2}{\sigma_w^2}I_N & \mathbf{k} \\ \mathbf{k}^{\mathsf{T}} & c + \frac{\sigma_\epsilon^2}{\sigma_w^2} \end{pmatrix})$$

# GP Regression

Manipulating Gaussians shows that

$$y_*|(y_1, \ldots, y_N) \sim \mathcal{N}(y_*|\mu, \boldsymbol{\Sigma})$$

where

$$\mu = \mathbf{k}^\intercal (K + \frac{\sigma_\epsilon^2}{\sigma_w^2} \mathbf{I}_N)^{-1} \mathbf{y}_o$$

$$\Sigma = \sigma_w^2 c + \sigma_\epsilon^2 - \sigma_w^2 \mathbf{k}^\intercal (K + \frac{\sigma_\epsilon^2}{\sigma_w^2} \mathbf{I}_N)^{-1} \mathbf{k}$$

# GP Regression

$$\mu = \mathbf{k}^{\mathsf{T}}(K + \frac{\sigma_\epsilon^2}{\sigma_w^2}\mathbf{I}_N)^{-1}\mathbf{y}_o$$

$\mu = \mathbf{k}^\mathsf{T}(K + \frac{\sigma_\epsilon^2}{\sigma_w^2}\mathbf{I}_N)^{-1}\mathbf{y}_o$

Some observations:

- Posterior mean depends on the ratio $\frac{\sigma_\epsilon^2}{\sigma_w^2}$

$\mu = \mathbf{k}^\intercal (K + \frac{\sigma_\epsilon^2}{\sigma_w^2} \mathbf{I}_N)^{-1} \mathbf{y}_o$

Some observations:

- Posterior mean depends on the ratio $\frac{\sigma_\epsilon^2}{\sigma_w^2}$
- Setting $\sigma_w^2 = 1$ and $\sigma_\epsilon^2 = \lambda$, we have KRR solution

$\mu = \mathbf{k}^{\mathsf{T}}(K + \frac{\sigma_\epsilon^2}{\sigma_w^2}\mathbf{I}_N)^{-1}\mathbf{y}_o$

Some observations:

- Posterior mean depends on the ratio $\frac{\sigma_\epsilon^2}{\sigma_w^2}$
- Setting $\sigma_w^2 = 1$ and $\sigma_\epsilon^2 = \lambda$, we have KRR solution
- $\mathrm{Cov}(y(\mathbf{x}), y(\mathbf{x}'))$ was in terms of inner products - can replace with any kernel function.

$\mu = \mathbf{k}^{\mathsf{T}}(K + \frac{\sigma_\epsilon^2}{\sigma_w^2}\mathbf{I}_N)^{-1}\mathbf{y}_o$

Some observations:

- Posterior mean depends on the ratio $\frac{\sigma_\epsilon^2}{\sigma_w^2}$
- Setting $\sigma_w^2 = 1$ and $\sigma_\epsilon^2 = \lambda$, we have KRR solution
- $\mathrm{Cov}(y(\mathbf{x}), y(\mathbf{x}'))$ was in terms of inner products - can replace with any kernel function.

Conclusion:

- KRR with kernel $k$ and regularisation $\lambda \subset$ GP regression with kernel $k' = k + \lambda\delta_{x=x'}$

$\mu = \mathbf{k}^{\mathsf{T}}(K + \frac{\sigma_\epsilon^2}{\sigma_w^2}\mathbf{I}_N)^{-1}\mathbf{y}_o$

Some observations:

- Posterior mean depends on the ratio $\frac{\sigma_\epsilon^2}{\sigma_w^2}$
- Setting $\sigma_w^2 = 1$ and $\sigma_\epsilon^2 = \lambda$, we have KRR solution
- $\text{Cov}(y(\mathbf{x}), y(\mathbf{x}'))$ was in terms of inner products - can replace with any kernel function.

Conclusion:

- KRR with kernel $k$ and regularisation $\lambda \subset$ GP regression with kernel $k' = k + \lambda\delta_{x=x'}$

In fact, if we use the kernel $k'$ for KRR without regularisation and just work through, we get the same answer[2].

---

[2]This is cheating really, because there is no unique optimum in this case

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$
- Want to infer a function $f : \mathcal{X} \longrightarrow \mathcal{Y}$ so that $f$ explains $\mathcal{D}$.

**Problem**:

- Given observations $\mathcal{D} = \{\mathbf{x}_i, y_i | i = 1, \ldots, N\}$ with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$
- Want to infer a function $f : \mathcal{X} \longrightarrow \mathcal{Y}$ so that $f$ explains $\mathcal{D}$.

**An approach**:

- Choose some set of candidate functions $\mathcal{F}$
- Choose some *loss function* $L(f, \mathcal{D})$ to penalise misfitting the data
- Choose some *complexity penalty* $\Omega(f)$ to prevent overfitting
- Find best $f \in \mathcal{F}$ to minimise sum:

$$\underset{f \in \mathcal{F}}{\arg\min}\, L(f, \mathcal{D}) + \Omega(f)$$

If $L(f, \mathcal{D}) = \sum_i L(f(x_i), y_i)$ then the problem is equivalent to

$$\arg\max_{f \in \mathcal{F}} \prod_i e^{-L(f(x_i), y_i)} e^{-\Omega(f)} \tag{*}$$

If $L(f, \mathcal{D}) = \sum_i L(f(x_i), y_i)$ then the problem is equivalent to

$$\arg\max_{f \in \mathcal{F}} \prod_i e^{-L(f(x_i), y_i)} e^{-\Omega(f)} \qquad (*)$$

If we can interpret

- $e^{-\Omega(f)}$ as a prior over $\mathcal{F}$
- $e^{-L(f(x_i), y_i)}$ as a likelihood

Then solving (*) is the same as performing MAP inference over $\mathcal{F}$.

# Kernel Ridge Regression as MAP

In Kernel Ridge Regression, the Representer theorem allowed us to restrict ourselves from $\mathcal{F} = \mathcal{H}$ to $\mathrm{span}\{k(\cdot, x_i)\}$. We parameterise $f$ by $\alpha$, and have $\Omega(f) = \lambda \alpha^{\mathsf{T}} K \alpha$. So we seek

$$\arg\max_{\alpha} \prod_i e^{-(f_\alpha(x_i) - y_i)^2} e^{-\lambda \alpha^{\mathsf{T}} K \alpha}$$

$$= \arg\max_{\alpha} \prod_i e^{-\frac{1}{2\lambda}(f_\alpha(x_i) - y_i)^2} e^{-\frac{1}{2}\alpha^{\mathsf{T}} K \alpha}$$

# Kernel Ridge Regression as MAP

In Kernel Ridge Regression, the Representer theorem allowed us to restrict ourselves from $\mathcal{F} = \mathcal{H}$ to $\text{span}\{k(\cdot, x_i)\}$. We parameterise $f$ by $\alpha$, and have $\Omega(f) = \lambda \alpha^\mathsf{T} K \alpha$. So we seek

$$\arg\max_\alpha \prod_i e^{-(f_\alpha(x_i) - y_i)^2} e^{-\lambda \alpha^\mathsf{T} K \alpha}$$

$$= \arg\max_\alpha \prod_i e^{-\frac{1}{2\lambda}(f_\alpha(x_i) - y_i)^2} e^{-\frac{1}{2}\alpha^\mathsf{T} K \alpha}$$

This is like finding the MAP solution in the model:

$$y | \alpha, x \sim \mathcal{N}(f_\alpha(x), \lambda) \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

$$y|\alpha, x \sim \mathcal{N}(f_\alpha(x), \lambda) \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

- Prior over $\alpha$ is Gaussian, $f_\alpha(x) = \mathbf{k}^\intercal \alpha \implies y$ is Gaussain
- $p(\alpha|\mathcal{D})$ also Gaussian due to self-conjugacy of Gaussian
- posterior over $y$ is Gaussian

$$y|\alpha, x \sim \mathcal{N}(f_\alpha(x), \lambda) \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

- Prior over $\alpha$ is Gaussian, $f_\alpha(x) = \mathbf{k}^\mathsf{T}\alpha \implies y$ is Gaussain
- $p(\alpha|\mathcal{D})$ also Gaussian due to self-conjugacy of Gaussian
- posterior over $y$ is Gaussian

So this model is a GP, and KRR gives its MAP solution

$$y|\alpha, x \sim \mathcal{N}(f_\alpha(x), \lambda) \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

- Prior over $\alpha$ is Gaussian, $f_\alpha(x) = \mathbf{k}^\mathsf{T}\alpha \implies y$ is Gaussain
- $p(\alpha|\mathcal{D})$ also Gaussian due to self-conjugacy of Gaussian
- posterior over $y$ is Gaussian

So this model is a GP, and KRR gives its MAP solution

Question: are there regression methods that are not strictly worse than GPs?

# Support Vector Regression

We can do the same as Kernel Ridge Regression but with a different loss function:

$$L(f(x), y) = \begin{cases} 0 & \text{if } |f(x) - y| < \epsilon \\ |f(x) - y| - \epsilon & \text{if } |f(x) - y| \geq \epsilon \end{cases}$$
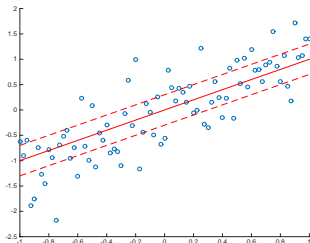
# Support Vector Regression

We can do the same as Kernel Ridge Regression but with a different loss function:

$$L(f(x), y) = \begin{cases} 0 & \text{if } |f(x) - y| < \epsilon \\ |f(x) - y| - \epsilon & \text{if } |f(x) - y| \geq \epsilon \end{cases}$$

**Why this might be a sensible $L$:**

1. Robust to outliers - linear rather than quadratic loss

2. Sparse solutions - any points inside $\epsilon$-tube around function are ignored

# Support Vector Regression

Want to solve:

$$\underset{f \in \mathcal{H}}{\arg\min} \sum_i L(f(x_i), y_i) + \lambda \|f\|_{\mathcal{H}}^2$$

# Support Vector Regression

Want to solve:

$$\underset{f \in \mathcal{H}}{\arg\min} \sum_i L(f(x_i), y_i) + \lambda \|f\|_{\mathcal{H}}^2$$

▶ Representer theorem $\implies$ solution lies in $\mathrm{span}\{k(\cdot, x_i)\}$
▶ Parameterise this subspace by $\alpha$ writing $f_\alpha(x) = \sum_i \alpha_i k(x, x_i)$

# Support Vector Regression

Want to solve:

$$\underset{f \in \mathcal{H}}{\arg\min} \sum_i L(f(x_i), y_i) + \lambda \|f\|_{\mathcal{H}}^2$$

- Representer theorem $\implies$ solution lies in $\mathrm{span}\{k(\cdot, x_i)\}$
- Parameterise this subspace by $\alpha$ writing $f_\alpha(x) = \sum_i \alpha_i k(x, x_i)$
- As before, problem reduces to:

$$\underset{\alpha}{\arg\min} \sum_i L(f_\alpha(x_i), y_i) + \lambda \alpha^\intercal K \alpha$$

# Support Vector Regression

This corresponds to the problem

$$\arg\max_{\alpha} \prod_i e^{-\frac{1}{2\lambda} L(f_\alpha(x_i), y_i)} e^{-\frac{1}{2}\alpha^\mathsf{T} K \alpha}$$

Equivalently, finding the MAP solution in the model:

$$p(y|\alpha, x) \propto e^{-\frac{1}{2\lambda} L(f_\alpha(x), y)} \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

This corresponds to the problem

$$\arg\max_{\alpha} \prod_i e^{-\frac{1}{2\lambda} L(f_\alpha(x_i), y_i)} e^{-\frac{1}{2}\alpha^\mathsf{T} K\alpha}$$

Equivalently, finding the MAP solution in the model:

$$p(y|\alpha, x) \propto e^{-\frac{1}{2\lambda} L(f_\alpha(x), y)} \qquad\qquad \alpha \sim \mathcal{N}(0, K^{-1})$$

- ▶ Prior on $\alpha$ is Gaussian, likelihood not Gaussian
- ▶ $\implies$ $y$ not Gaussian, posterior $p(\alpha|\mathcal{D})$ not Gaussian
- ▶ $\implies$ latent function values $f_\alpha(x) = \mathbf{k}^\mathsf{T}\alpha$ will not be Gaussian

So Support Vector Regression is distinct from Gaussian Process Regression.

1. Derived Kernel Ridge Regression

# Conclusion

1. Derived Kernel Ridge Regression
2. KRR solution same as GP posterior mean (so KRR $\subset$ GP)

# Conclusion

1. Derived Kernel Ridge Regression
2. KRR solution same as GP posterior mean (so KRR $\subset$ GP)
3. KRR is like MAP inference in GP model

# Conclusion

1. Derived Kernel Ridge Regression
2. KRR solution same as GP posterior mean (so KRR $\subset$ GP)
3. KRR is like MAP inference in GP model
4. Support Vector Regression is not comparable to GP regression

UNIVERSITY OF
CAMBRIDGE