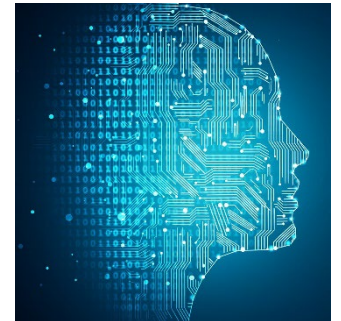


Machine Learning

Kernel Ridge Regression



Kevin Moon (kevin.moon@usu.edu)
STAT/CS 5810/6655



Outline



1. Review of kernels
2. Kernel ridge regression
3. Kernel nearest centroid classifier
4. Kernel ridge regression revisited

Nonlinear Feature Maps

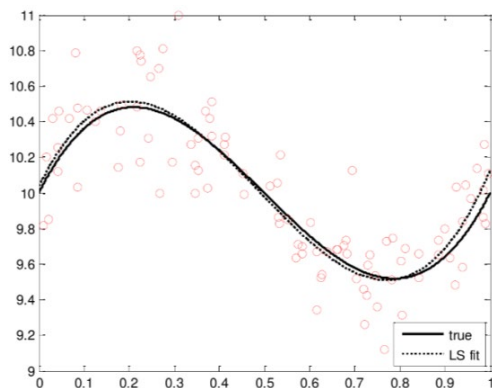


- One way to create a nonlinear method for regression or classification is to first transform the feature vector via a *nonlinear feature map*

$$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

and apply a linear method to the transformed data $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$.

- **Examples:** polynomial regression function, circular classifier
- Problem with the above approach: m can explode as d increases
- This makes it prohibitive to compute/store/manipulate Φ directly
- This is where *inner product kernels* save the day



Inner Product Kernels



- We say $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is an *inner product kernel* if \exists an inner product space V and a feature map $\Phi : \mathbb{R}^d \rightarrow V$ such that

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$$

- Homogeneous polynomial kernel

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^p$$

$$\Phi(\mathbf{u}) = \text{monomials in } u^{(1)}, \dots, u^{(d)} \text{ of degree } p$$

- Inhomogeneous polynomial

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + c)^p, \quad c > 0$$

$$\Phi(\mathbf{u}) = \text{monomials in } u^{(1)}, \dots, u^{(d)} \text{ up to degree } p$$

- Gaussian kernel

$$k(\mathbf{u}, \mathbf{v}) = \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{u} - \mathbf{v}\|^2 \right), \quad \sigma > 0, \quad m = \infty$$

The Kernel Trick



- A machine learning algorithm is said to be *kernelizable* if it can be formulated such that feature vectors (i.e., the training data and an arbitrary test instance) only appear via inner products $\langle \mathbf{x}, \mathbf{x}' \rangle$ with other feature vectors.
- Suppose Φ is a feature map associated with an inner product kernel k .
- If we apply a kernelizable algorithm to the training data $(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_n), y_n)$, then we can formulate the algorithm such that transformed features only appear via inner products $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$ with other transformed features
- Can implement by evaluating $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$ which eliminates the need to ever compute $\Phi(\mathbf{x})$ explicitly
- In practice, we don't even need to know Φ , we just need to specify an inner product kernel k

Kernel Ridge Regression (1)



- Ridge regression solves

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x} - b)^2 + \lambda \|\mathbf{w}\|^2$$

- We will show that ridge regression is kernelizable and use the kernel trick to extend it to a nonlinear regression method called *kernel ridge regression*
- To simplify the presentation, we'll first consider ridge regression without offset

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x})^2 + \lambda \|\mathbf{w}\|^2$$

- Since some Φ contain a constant term, as with the inhomogeneous polynomial, an offset isn't always needed

Kernel Ridge Regression (2)



- We can rewrite the objective function as follows (without changing the solution):

$$\|\mathbf{y} - X\mathbf{w}\|^2 + n\lambda\|\mathbf{w}\|^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad X = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(d)} \\ \vdots & & \vdots \\ x_n^{(1)} & \dots & x_n^{(d)} \end{bmatrix}$$

Kernel Ridge Regression (3)



- The solution is

$$\hat{\mathbf{w}} = (X^T X + n\lambda I)^{-1} X^T \mathbf{y}$$

- Problem: $[X^T X]_{ij} \neq \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
 - I.e. $X^T X$ is not the gram matrix of the training data
- Idea: apply the matrix inversion lemma:

$$(P + QRS)^{-1} = P^{-1} - P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}SP^{-1}$$

with

$$\begin{aligned} P &= n\lambda I = \mu I, & Q &= X^T, \\ S &= X, & R &= I \end{aligned}$$

Kernel Ridge Regression (4)



- Plugging in:

$$\begin{aligned}(\mu I + X^T X)^{-1} &= \frac{1}{\mu} I - \frac{1}{\mu} I X^T \left(I + \frac{1}{\mu} X X^T \right)^{-1} X \frac{1}{\mu} I \\ &= \frac{1}{\mu} \left(I - X^T (\mu I + X X^T)^{-1} X \right)\end{aligned}$$

- Therefore,

$$\begin{aligned}\hat{\mathbf{w}} &= (X^T X + n\lambda I)^{-1} X^T \mathbf{y} = \frac{1}{\mu} \left(X^T - X^T (\mu I + X X^T)^{-1} X X^T \right) \mathbf{y} \\ &= \frac{1}{\mu} \left(X^T - X^T (G + \mu I)^{-1} G \right) \mathbf{y}\end{aligned}$$

where

$$G = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \dots & \langle \mathbf{x}_1, \mathbf{x}_n \rangle \\ \vdots & & \vdots \\ \langle \mathbf{x}_n, \mathbf{x}_1 \rangle & \dots & \langle \mathbf{x}_n, \mathbf{x}_n \rangle \end{bmatrix}$$

Kernel Ridge Regression (5)



- The ridge regression function is therefore

$$\begin{aligned}\hat{f}(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} \\ &= \frac{1}{\mu} \mathbf{y}^T (X - G(G + \mu I)^{-1} X) \mathbf{x} \\ &= \frac{1}{\mu} \mathbf{y}^T (I - G(G + \mu I)^{-1}) X \mathbf{x}\end{aligned}$$

Kernel Ridge Regression (6)



- Simplifying $\hat{f}(\mathbf{x}) = \frac{1}{\mu} \mathbf{y}^T (I - G(G + \mu I)^{-1}) X \mathbf{x}$

$$\begin{aligned} I - G(G + \mu I)^{-1} &= (G + \mu I)(G + \mu I)^{-1} - G(G + \mu I)^{-1} \\ &= (G + \mu I - G)(G + \mu I)^{-1} \\ &= \mu(G + \mu I)^{-1} \end{aligned}$$

$$X \mathbf{x} = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{x}_n, \mathbf{x} \rangle \end{bmatrix}$$

Kernel Ridge Regression (7)



- In summary, we have expressed the solution of ridge regression w/o offset as

$$\hat{f}(\mathbf{x}) = \mathbf{y}^T (G + n\lambda I)^{-1} \mathbf{g}(\mathbf{x})$$

where

$$XX^T = G = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \dots & \langle \mathbf{x}_1, \mathbf{x}_n \rangle \\ \vdots & & \vdots \\ \langle \mathbf{x}_n, \mathbf{x}_1 \rangle & \dots & \langle \mathbf{x}_n, \mathbf{x}_n \rangle \end{bmatrix}$$

$$\mathbf{g}(\mathbf{x}) = X\mathbf{x} = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{x}_n, \mathbf{x} \rangle \end{bmatrix}$$

Kernel Ridge Regression (8)



- Now suppose we first transform the feature vectors by a nonlinear feature map Φ , and then apply ridge regression (w/o offset) in the new feature space.
- Our regression function estimate is then

$$\hat{f}(\mathbf{x}) = \mathbf{y}^T (K + n\lambda I)^{-1} \mathbf{k}(\mathbf{x})$$

where

$$K = \begin{bmatrix} \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_1) \rangle & \dots & \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_n) \rangle \\ \vdots & & \vdots \\ \langle \Phi(\mathbf{x}_n), \Phi(\mathbf{x}_1) \rangle & \dots & \langle \Phi(\mathbf{x}_n), \Phi(\mathbf{x}_n) \rangle \end{bmatrix}$$

$$\mathbf{k}(\mathbf{x}) = \begin{bmatrix} \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}) \rangle \\ \vdots \\ \langle \Phi(\mathbf{x}_n), \Phi(\mathbf{x}) \rangle \end{bmatrix}$$

Kernel Ridge Regression (9)



- If Φ is the feature map associated with an inner product kernel k , then

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}, \quad k(\mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}) \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}) \end{bmatrix}$$

- Conclusion: we don't need to compute $\Phi(\mathbf{x})$

KRR (w/o offset) + Gaussian Kernel

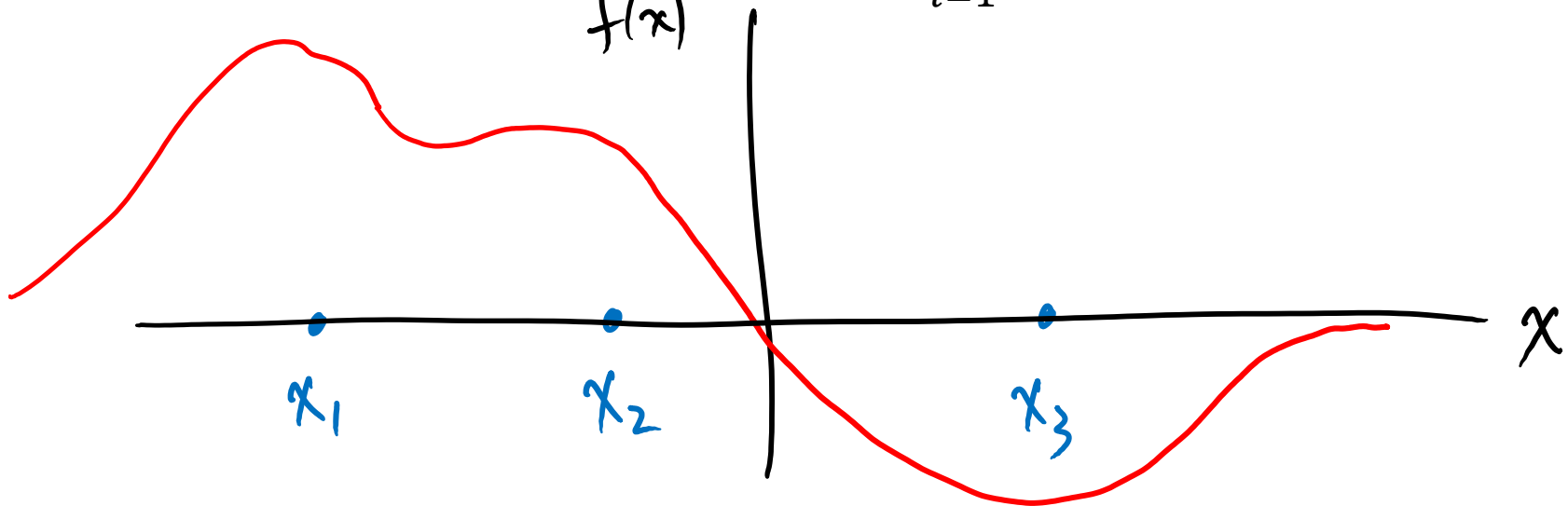


- Consider the Gaussian kernel

$$k(u, v) = \exp\left(-\frac{1}{2\sigma^2} \|u - v\|^2\right), \sigma > 0$$

- The KRR (w/o offset) predictor can be expressed as

$$\hat{f}(x) = \alpha^T k(x) = \sum_{i=1}^n \alpha_i k(x_i, x)$$



- Although the Φ associated with the Gaussian kernel does not contain a constant, KRR with the Gaussian kernel often performs well w/o the offset

Group Exercise



Nearest centroid classifier: Given training data $(\mathbf{x}_i, y_i), i = 1, \dots, n$, $y_i \in \{-1, +1\}$, define the *centroids*

$$\mathbf{m}_- = \frac{1}{n_-} \sum_{i: y_i = -1} \mathbf{x}_i \quad \mathbf{m}_+ = \frac{1}{n_+} \sum_{i: y_i = +1} \mathbf{x}_i$$

where $n_- = |\{i : y_i = -1\}|$ and $n_+ = |\{i : y_i = +1\}|$. The nearest centroid classifier is

$$\begin{aligned} f(\mathbf{x}) &= \text{sign}\{\|\mathbf{x} - \mathbf{m}_-\|^2 - \|\mathbf{x} - \mathbf{m}_+\|^2\} \\ &= \text{sign} \left\{ \left(\mathbf{x} - \frac{\mathbf{m}_+ + \mathbf{m}_-}{2} \right)^T (\mathbf{m}_+ - \mathbf{m}_-) \right\} \end{aligned}$$

1. Kernelize the nearest centroid classifier.

KRR with Offset



- Now assume there is an offset:

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|^2$$

- Recall the solution is

$$\begin{aligned}\hat{\mathbf{w}} &= (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + n\lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}} \\ \hat{b} &= \bar{y} - \hat{\mathbf{w}}^T \bar{\mathbf{x}}\end{aligned}$$

where

$$\begin{aligned}\tilde{\mathbf{y}} &= \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix}, & \tilde{\mathbf{X}} &= \begin{bmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_n^T \end{bmatrix}, & \tilde{\mathbf{x}}_i &= \mathbf{x}_i - \bar{\mathbf{x}} \in \mathbb{R}^d \\ & & & & \tilde{y}_i &= y_i - \bar{y} \in \mathbb{R} \\ & & & & \bar{\mathbf{x}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \\ & & & & \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i\end{aligned}$$

KRR with Offset



- Using the exact same algebra as before, it can be shown that

$$\hat{f}(\mathbf{x}) = \bar{y} + \tilde{\mathbf{y}}^T (\tilde{G} + n\lambda I)^{-1} \tilde{\mathbf{g}}(\mathbf{x})$$

where

$$\tilde{G} = \begin{bmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_1 \rangle & \dots & \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_n \rangle \\ \vdots & & \vdots \\ \langle \tilde{\mathbf{x}}_n, \tilde{\mathbf{x}}_1 \rangle & \dots & \langle \tilde{\mathbf{x}}_n, \tilde{\mathbf{x}}_n \rangle \end{bmatrix}, \quad \tilde{\mathbf{g}}(\mathbf{x}) = \begin{bmatrix} \langle \tilde{\mathbf{x}}_1, \mathbf{x} - \bar{\mathbf{x}} \rangle \\ \vdots \\ \langle \tilde{\mathbf{x}}_n, \mathbf{x} - \bar{\mathbf{x}} \rangle \end{bmatrix}$$

- Now suppose we first transform the data with Φ
- It's tempting to substitute $\langle \tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \rangle$ with $\langle \Phi(\tilde{\mathbf{x}}), \Phi(\tilde{\mathbf{x}}') \rangle$

KRR with Offset



- To kernelize, observe

$$\begin{aligned}\langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle &= \langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{x}_j - \bar{\mathbf{x}} \rangle \\ &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \frac{1}{n} \sum_{r=1}^n \langle \mathbf{x}_i, \mathbf{x}_r \rangle - \frac{1}{n} \sum_{s=1}^n \langle \mathbf{x}_s, \mathbf{x}_j \rangle + \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n \langle \mathbf{x}_r, \mathbf{x}_s \rangle\end{aligned}$$

And

$$\begin{aligned}\langle \tilde{\mathbf{x}}_i, \mathbf{x} - \bar{\mathbf{x}} \rangle &= \langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{x} - \bar{\mathbf{x}} \rangle \\ &= \langle \mathbf{x}_i, \mathbf{x} \rangle - \frac{1}{n} \sum_{r=1}^n \langle \mathbf{x}_i, \mathbf{x}_r \rangle - \frac{1}{n} \sum_{s=1}^n \langle \mathbf{x}, \mathbf{x}_s \rangle + \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n \langle \mathbf{x}_r, \mathbf{x}_s \rangle\end{aligned}$$

- Now replace $\langle \mathbf{x}, \mathbf{x}' \rangle$ with $k(\mathbf{x}, \mathbf{x}')$ throughout the above expressions

KRR with Offset: alternative View



- Define the “mean-centered feature map”

$$\tilde{\Phi}(\mathbf{x}) := \Phi(\mathbf{x}) - \frac{1}{n} \sum_{\ell=1}^n \Phi(\mathbf{x}_{\ell})$$

- KRR with offset can equivalently be expressed

$$\hat{f}(\mathbf{x}) = \bar{y} + \tilde{\mathbf{y}}^T (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \tilde{\mathbf{k}}(\mathbf{x})$$

where the (i, j) entry of $\tilde{\mathbf{K}}$ is

$$\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}_j) \rangle$$

and the i -th entry of $\tilde{\mathbf{k}}(\mathbf{x})$ is

$$\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}) \rangle$$

- These values are computed using the formulas on the previous slide, after replacing dot products with kernels.