

Foundations of Machine Learning

AI2000 and AI5000

FoML-31

Kernelized Linear Models

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

- Intro to ML and Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Supervised learning
 - a. Linear Regression with basis functions
 - b. Bias-Variance Decomposition
 - c. Decision Theory & three broad classification strategies
 - d. Neural Networks
- Unsupervised learning
 - a. K-Means, Hierarchical, GMM for clustering, and PCA



For today

- Equivalent Kernel
- Kernelizing Linear models



Recap

- Bayesian Regression (foml-11)



Equivalent Kernel



Equivalent Kernel formulation

- The predictive distribution: $p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}.\end{aligned}$$

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x})$$



Equivalent Kernel formulation

- The predictive distribution: $p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}.\end{aligned}\quad \sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x})$$

- predictive mean: $y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}) = \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} = \sum_{n=1}^N \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_n) t_n$

$$k(x, x_n) =$$

$$y(\mathbf{x}, \mathbf{m}_N) =$$



Equivalent Kernel formulation

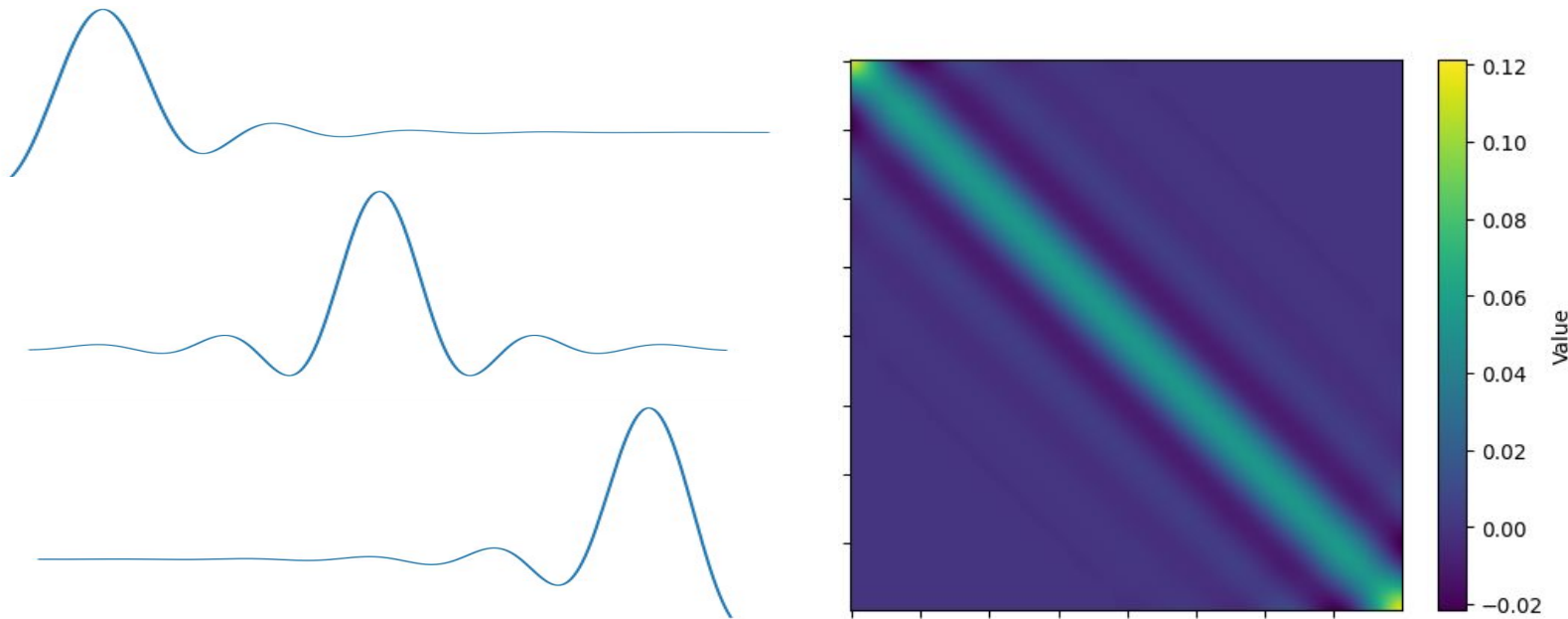
$$k(\mathbf{x}, \mathbf{x}_n) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n)$$

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

- K - smoother matrix or equivalent kernel
 - It depends on all the input samples
- Functions that make predictions by taking linear combinations of the training set target values - linear smoothers



Equivalent kernel for Gaussian Basis functions



Training samples close to x contribute more!



Covariance between two predictions

$$\text{Cov}[t_1, t_2 | \mathbf{x}_1, \mathbf{x}_2] =$$



Alternate approach to parametric modeling

- Instead of working with basis functions
- Define a localized kernel to make predictions for new points
 - Gaussian Processes

Summary - Parametric Models

- Use fixed basis function to project the data
 - Learning: regression, classification
- Learnable basis functions: neural networks
- Training
 - MLE, MAP \rightarrow point estimate W
 - Full Bayesian \rightarrow posterior on W
- Test time
 - Don't need the training data
 - Work with W or its distribution



Memory-based Methods

- Training data is kept and used for inference
 - KDE
 - KNN
- Fast 'training', slow 'inference'



Non-parametric Kernel Methods



Non-parametric methods

- Kernel methods
 - Use training data for test time predictions
- 'Dual representation' for the linear parametric models
 - Equivalent kernels

Kernelized Ridge Regression

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Solution to \mathbf{w} takes a form of linear combination of basis vectors



Kernelized Ridge Regression

- Instead of working with \mathbf{w} , let us work with \mathbf{a}

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad \mathbf{w} = \Phi^T \mathbf{a}$$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$



Kernelized Ridge Regression

- Instead of working with w , let us work with a

$$J(a) = \frac{1}{2}a^T \Phi \Phi^T \Phi \Phi^T a - a^T \Phi \Phi^T t + \frac{1}{2}t^T t + \frac{\lambda}{2}a^T \Phi \Phi^T a$$

- Introduce a gram matrix K $K_{nm} = \phi(x_n)^T \phi(x_m) = k(x_n, x_m)$



Kernelized Ridge Regression

- Optimizing for \mathbf{a}

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$$

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.$$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$



Primal and Dual perspectives

$$\mathbf{w} = \Phi^T \mathbf{a}$$

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.$$

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$y(\mathbf{x}, \mathbf{a}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x})$$

Primal and Dual perspectives

$$\mathbf{w} = \Phi^T \mathbf{a}$$

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.$$

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$y(\mathbf{x}, \mathbf{a}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x})$$

- Compute (train)
 - $O(M^3)$ vs. $O(N^3)$
- Compute (inference)
 - $O(M)$ vs. $O(NM)$



Primal and Dual perspectives



Next

- Kernel methods with 'sparsity' in solutions

