

Foundations of Machine Learning

AI2000 and AI5000

FoML-31
Kernelized Linear Models

Dr. Konda Reddy Mopuri
Department of AI, IIT Hyderabad
July-Nov 2025

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



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રેરણાખાડ
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ
Indian Institute of Technology Hyderabad

Recap

- Bayesian Regression (foml-11)

Equivalent Kernel



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રેરણાખાર્ડ
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ
Indian Institute of Technology Hyderabad



Equivalent Kernel formulation

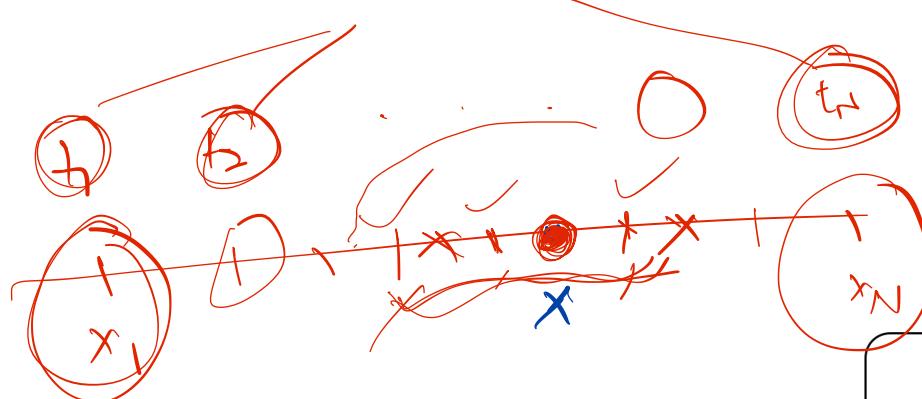
- The predictive distribution: $p(t|\bar{\mathbf{x}}, \bar{\mathbf{t}}, \bar{\alpha}, \bar{\beta}) = \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$

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

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

$K(x, x_n)$

$$f(x) = k$$



$$\Phi = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_N) \end{bmatrix}$$

Equivalent Kernel formulation

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

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

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

- predictive mean:

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

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

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



Equivalent Kernel formulation

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

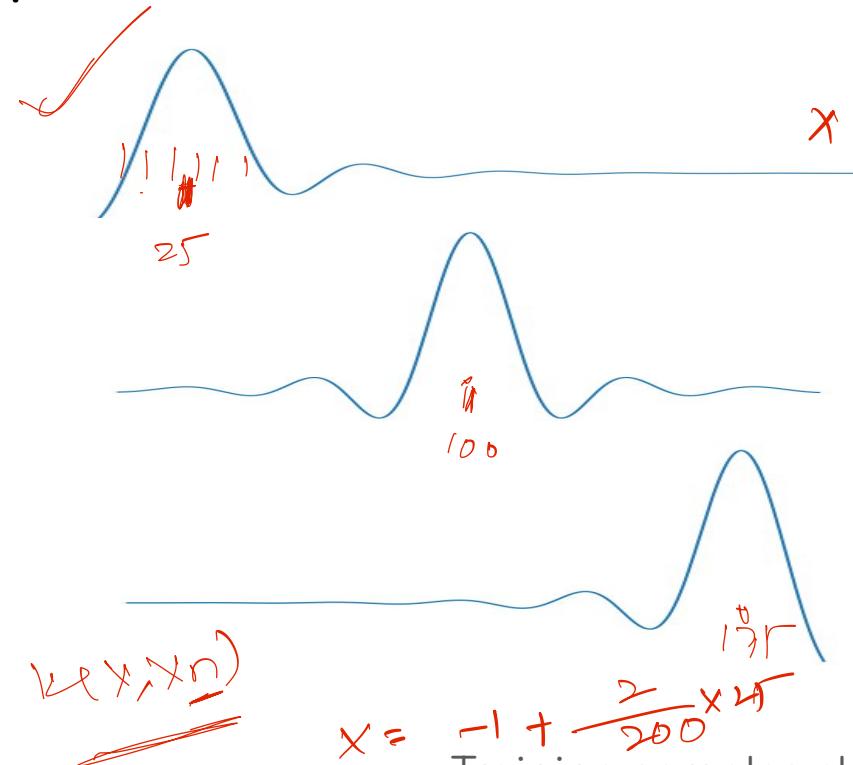
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

- 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

$$\Phi \left(\begin{bmatrix} \gamma_2 \\ \vdots \\ \gamma_N \\ \phi(\mathbf{x}) \end{bmatrix} \right)$$

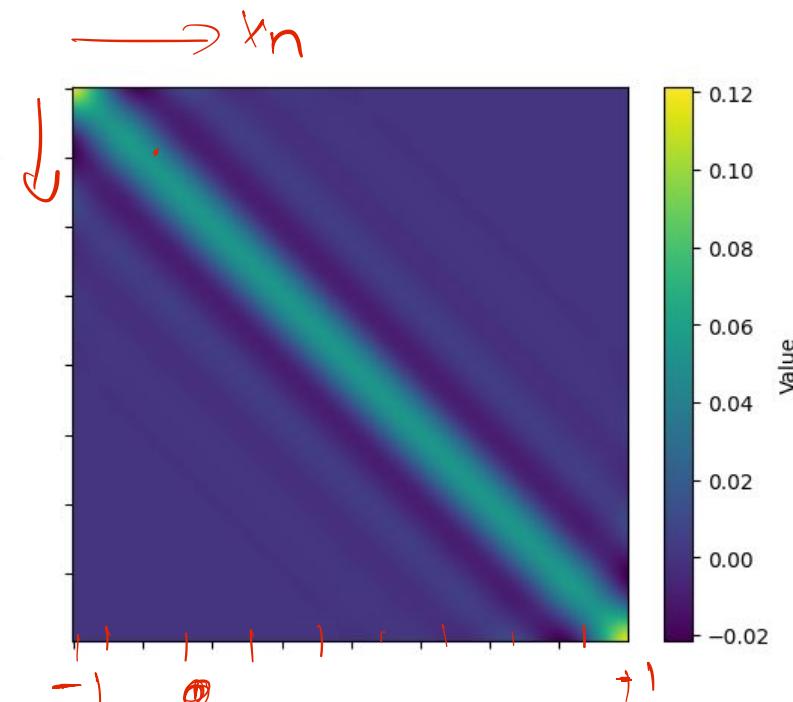


Equivalent kernel for Gaussian Basis functions



$$n < 100$$

$$200 \quad [-1, 1]$$



Training samples close to x contribute more!



Covariance between two predictions

$$\begin{aligned}
 \text{Cov}[t_1, t_2 | \mathbf{x}_1, \mathbf{x}_2] &= \text{Cov}_{\mathbf{w}} \left(\underbrace{\phi(\mathbf{x}_1)^\top \mathbf{w}}_{\sim N(\mathbf{0}, \Sigma)}, \underbrace{\mathbf{w}^\top \phi(\mathbf{x}_2)}_{\sim N(\mathbf{0}, \Sigma)} \right) \\
 &= E_{\mathbf{w}} \left[\underbrace{\phi(\mathbf{x}_1)^\top \mathbf{w} \mathbf{w}^\top \phi(\mathbf{x}_2)}_{\text{cov}} \right] - E_{\mathbf{w}}[\phi(\mathbf{x}_1)^\top \mathbf{w}] E_{\mathbf{w}}[\mathbf{w}^\top \phi(\mathbf{x}_2)] \\
 &= \phi(\mathbf{x}_1)^\top \underbrace{E[\mathbf{w} \mathbf{w}^\top]}_{\text{cov}} \phi(\mathbf{x}_2) - \phi(\mathbf{x}_1)^\top E[\mathbf{w}] E[\mathbf{w}^\top] \phi(\mathbf{x}_2) \\
 &= \phi(\mathbf{x}_1)^\top \left(E[\mathbf{w} \mathbf{w}^\top] - E[\mathbf{w}] E[\mathbf{w}^\top] \right) \phi(\mathbf{x}_2) \\
 &= \phi(\mathbf{x}_1)^\top \underbrace{K(\mathbf{x}_1, \mathbf{x}_2)}_{\phi(\mathbf{x}_1)^\top K_N \phi(\mathbf{x}_2) = \beta K(\mathbf{x}_1, \mathbf{x}_2)}
 \end{aligned}$$



भारतीय नोर्मेटिक विज्ञान संस्था हैदराबाद
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



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 → point estimate W
 - Full Bayesian → posterior on W
- Test time
 - Don't need the training data
 - Work with W or its distribution



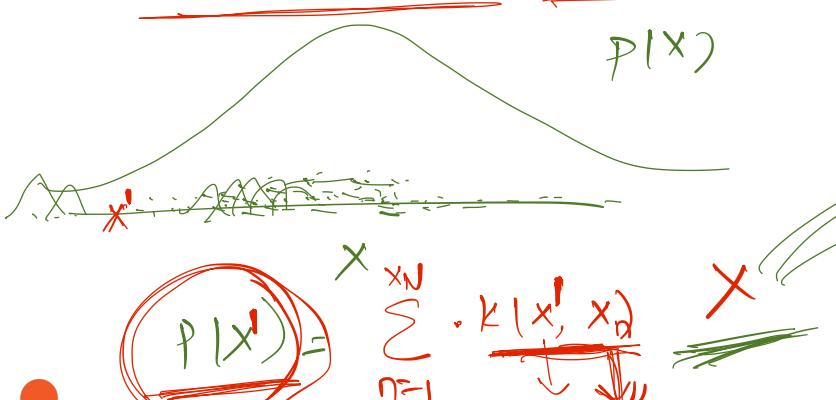
भारतीय नॉर्किंग विज्ञान संस्था पैदार्थाभार्द

भारतीय प्रौद्योगिकी संस्थान हैदराबाद

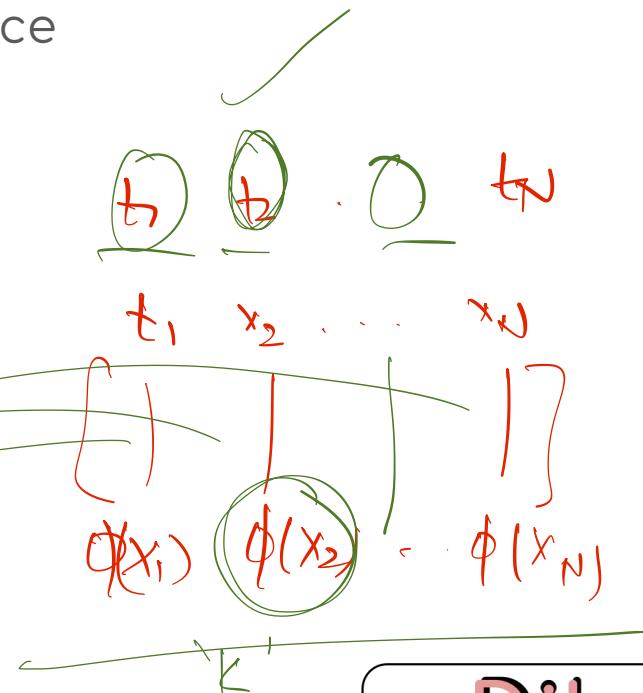
Indian Institute of Technology Hyderabad

Memory-based Methods

- Training data is kept and used for inference
 - KDE
 - KNN
- Fast ‘training’, slow ‘inference’
 - Fast ‘training’
 - Slow ‘inference’



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



Non-parametric Kernel Methods



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રેરણાખાડ
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ
Indian Institute of Technology Hyderabad



Non-parametric methods

- Kernel methods ✓
- Use training data for test time predictions
- ‘Dual representation’ for the linear parametric models
○ Equivalent kernels



Kernelized Ridge Regression

$\underset{\mathbf{w}}{\operatorname{arg\,min}}$

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

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

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) = 0 \quad \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) + \lambda \mathbf{w} = 0$$

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

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

$\Phi^T \in \mathbb{R}^{N \times N}$ $\mathbf{a} \in \mathbb{R}^N$



Kernelized Ridge Regression

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

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

\checkmark

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

$w \rightarrow a$
 $N \times 1$ $N \times 1$



Kernelized Ridge Regression

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

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

- Introduce a gram matrix \underline{K}

$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$



$$\begin{matrix} K = \phi \phi^T \\ \hline N \times N & N \times M & M \times N \end{matrix}$$



Kernelized Ridge Regression

- Optimizing for \underline{a}

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

$$\frac{\partial}{\partial \underline{a}} J(\underline{a})$$

$$K \underline{a} - K \underline{t} + \lambda K \underline{a} = 0$$

$$K \underline{t} = K \underline{a} + \lambda K \underline{a}$$

$$\underline{t} = \frac{K \underline{a} + \lambda I_N \underline{a}}{(K + \lambda I_N)^{-1}}$$

$$\underline{a} = (K + \lambda I_N)^{-1} \underline{t}$$

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



Primal and Dual perspectives

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

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

Inference

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

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

Inform



Primal and Dual perspectives

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

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

$$\left\{ \begin{array}{l} \mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \\ \text{---} \\ M \times M \end{array} \right.$$

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

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

Primal is
Dual is
Costly

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

$\frac{N \gg M}{\# \text{data} \quad \# \text{dim}}$



Primal and Dual perspectives

-

Next

- Kernel methods with 'sparsity' in solutions



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રેરણાખાડ
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ
Indian Institute of Technology Hyderabad