

Kernel Based Regression

- Generalized notion of predictor variables
- Use ridge on generalized variables
- Kernels useful in extending classical regression tools to nontraditional data.
 - pictures - graphs

Goal: Find $\hat{h}: V \rightarrow \mathbb{R}$ describing how x vary across V .

Def: $K: V \times V \rightarrow \mathbb{R}$ is a psd kernel iff for all submatrices
is (a) symmetric (b) psd

Remark: Kernels create similarity matrices

- Assume our function comes from the class

$$\mathcal{H}_K = \left\{ h \in \mathbb{R}^{N_V} : \underbrace{h = \mathbb{I}\beta}_{\text{linearity}} \text{ and } \underbrace{\beta^T \Delta^{-1} \beta}_{\text{penalty}} < \infty \right\}$$

$K = \mathbb{I} \Delta \mathbb{I}^T$ - build \mathcal{H}_K from eigen decomp of K .

Remark: Choose K which induces class of functions \mathcal{H}_K

Solve this problem via

$$\min_{\beta} \left[\sum_{i \in \text{data}} C(x_i, (\mathbb{I}\beta)_i) + \lambda \beta^T \Delta^{-1} \beta \right]$$

Representer Thm: (Don't need to evaluate K .) \hat{h} takes the form $h = K_{\text{obs.}}^{(N_{\text{tr}}, n)} \alpha$

where $K^{(N_{\text{tr}}, n)}$ is an evaluation on $V \times V^{\text{obs}}$

Reduces the optimization

$$\min_{\alpha} \left[\sum_{i \in \text{obs}} C(x_i, (K^{(n)} \alpha)_i) + \lambda \alpha^T K^{(n)} \alpha \right]$$

where $K^{(n)}$ is $n \times n$

Under squared error loss

$$\min_{\alpha} \left[(x - K^{(n)} \alpha)^T (x - K^{(n)} \alpha) + \lambda \alpha^T K^{(n)} \alpha \right]$$

Change of variables $m = \mathbb{E}_n \Delta_n^{-1/2}$ $\theta = \Delta_n^{-1/2} \mathbb{E}_n^T \alpha$ for $K^n = \mathbb{E}_n \Delta_n \mathbb{E}_n^T$

$$\min_{\theta} \left[(x - m\theta)^T (x - m\theta) + \lambda \theta^T \theta \right]$$

$$\hat{\theta}_{\text{RR}} = (m^T m + \lambda I)^{-1} m^T x^{\text{obs.}}$$

$$\rightarrow \hat{\alpha} = \mathbb{E}_n \Delta_n^{-1/2} \hat{\theta}_{\text{RR}}$$

and our final solution is $\boxed{\hat{h} = K^{(N_{\text{tr}}, n)} \hat{\alpha}}$

In designing a kernel: - Captures similarity among vertices

- p.s.d.

Ex: Laplacian Kernel: $K = L$

under this kernel, the penalty β has the form

$$\beta^T \Delta^{-1} \beta = \beta^T \mathbb{I}^T \mathbb{I} \Delta^{-1} \mathbb{I} \beta \\ = h^T K h$$

$$= h^T L h$$

$$= \sum_{(i,j) \in E} (h_i - h_j)^2$$

So this kernel gives a smooth penalty over the graph.

A large class of kernel $K = \sum_{i=1}^N r_i(x_i) \phi_i \phi_i^T$ based on this