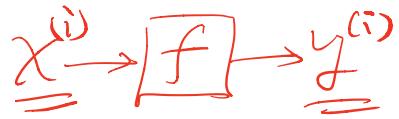


Nonlinear Function Approximation: Kernel Method

Qiang Liu
UT Austin

Supervised Learning Framework

- Given training data $\mathcal{D} = \{\underline{x}_i, \underline{y}_i\}_{i=1}^n$.
- Want to find $f(\underline{x})$, such that $\underline{y}_i \approx f(\underline{x}_i)$.
- Empirical risk minimization:
 - Decide a function class \mathcal{F} .
 - Define an empirical loss function $L(f; \mathcal{D})$ for $f \in \mathcal{F}$
 - Solve optimization:
$$\hat{f} = \arg \min_{f \in \mathcal{F}} L(f; \mathcal{D}).$$



$$\mathcal{F} \triangleq \left\{ f_{\theta}(x) = \sum_{\ell=1}^d \theta_{\ell} x_{\ell} + \theta_0 \mid \forall \theta_{\ell} \in \mathbb{R} \right\}$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\min_{\theta} L(f_{\theta}, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - f_{\theta}(x^{(i)}))^2$$

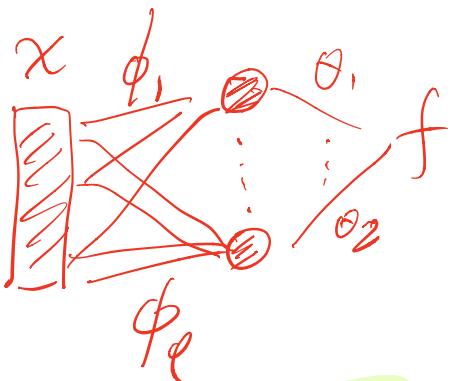
$$\text{Linear: } f_{\theta}(x) = \sum_{\ell=1}^d \theta_{\ell} \cancel{x_{\ell}}$$

$$\text{Nonlinear: } f_{\theta}(x) = \sum_{\ell=1}^{d'} \theta_{\ell} \cancel{\phi_{\ell}(x)}$$

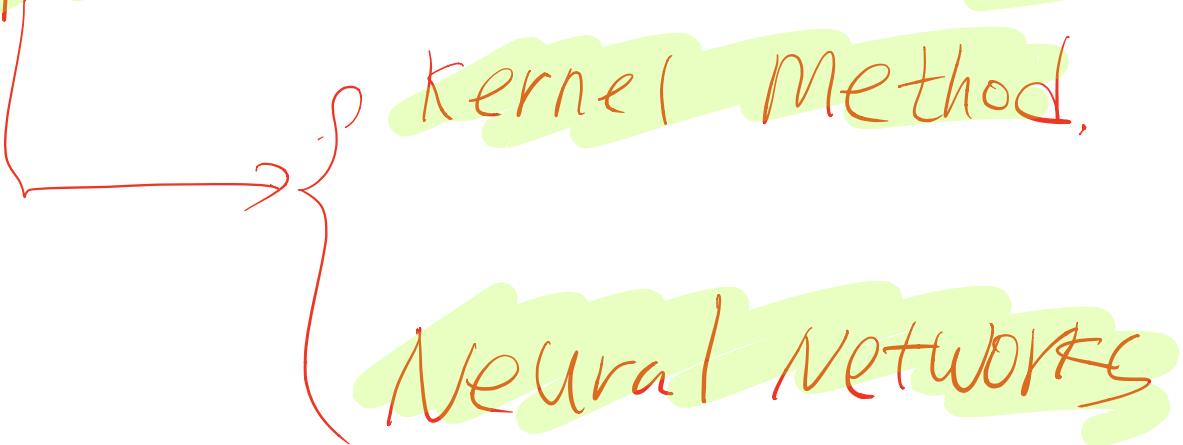
ϕ_{ℓ} : Basis function

~~Ch. Basis Function~~

$$\underline{\phi}_l(x) = x^l \Rightarrow 1, x, x^2, \dots$$



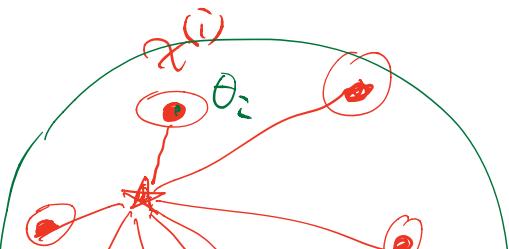
Adaptive Basis function



Kernel Method

$$D = \{x^{(i)}, y^{(i)}\}$$

Similarity function



$$\underline{K(x, x')} : X \times X \mapsto \mathbb{R}$$

Example:

$$K(x, x') = \exp\left(-\frac{1}{2h^2} \|x - x'\|^2\right)$$

Gaussian Radial Basis function (RBF)

kernel. h : Bandwidth

Kernel = A two variable that is
symmetric

$$K(x, x') = K(x', x)$$

$$K(x, x') = \exp\left(-\frac{1}{h^2} \|x - x'\|^2\right)$$

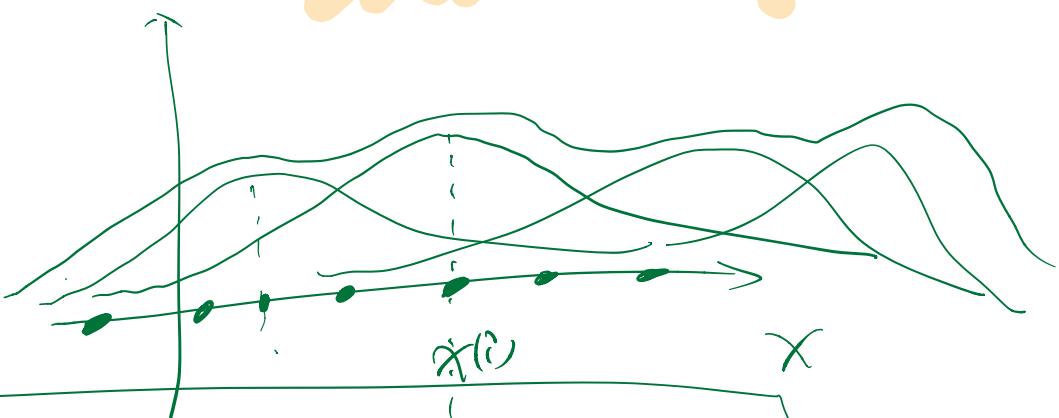
Laplace Kernel.

$$\phi(x) = \begin{bmatrix} K(x, x_1) \\ K(x, x_2) \\ \vdots \end{bmatrix} \quad \Theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\left[\begin{array}{c} k(x, x_1) \\ \vdots \\ k(x, x_n) \end{array} \right] \quad \left\{ \theta_n \right\}$$

$$f_\theta(x) = \underline{\theta}^T \underline{\phi}(x) = \sum_{i=1}^n \underline{\theta}_i \underline{\phi}_i(x)$$

$$= \sum_{i=1}^n \underline{\theta}_i k(x, x_i)$$



\mathcal{F} : adaptive with data D .

$$\dim(\mathcal{F}) = n \rightarrow +\infty$$

Nonparametric Method.



$$\min_{\theta} L(\theta) = \sum_{j=1}^n (y^{(j)} - \sum_{i=1}^n \theta_i k(x^{(j)}, x^{(i)}))$$

$$= \|Y - K\theta\|_2^2$$

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}, \quad K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix}_{n \times n}$$

$$\Rightarrow \nabla L(\theta) = 2K(K\theta - Y)$$

$$\Rightarrow K^T K \theta = K^T Y$$

$$\Rightarrow \boxed{\theta = K^{-1} Y} \quad (\text{Assume } K \text{ is invertible})$$

$$\Rightarrow \min_{\theta} \|Y - K\theta\|_2^2 + \alpha \underline{\Phi(\theta)}$$

$L_2: \|\theta\|_2^2$

$$= \theta^T \theta$$

$\|\theta\|_K^2 = \theta^T K \theta$

why $\|\theta\|_K^2$ vs. $\|\theta\|_2^2$?

$$D = \{ \underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)} \}$$

Assume $\underline{x}^{(1)} = \underline{x}^{(2)} \neq \underline{x}^{(3)}$

$x^{(1)} = x^{(2)}$

$$f(x) = \theta_1 k(x, x^{(1)}) + \theta_2 k(x, x^{(2)}) + \theta_3 k(x, x^{(3)})$$

$$= \underline{(\theta_1 + \theta_2)} \underline{K(x, x^{(1)})} + \theta_3 K(x, x^{(3)})$$

$$\Rightarrow \|\theta\|_2^2 = \underline{\theta_1^2} + \underline{\theta_2^2} + \underline{\theta_3^2}$$

Better: $\underline{(\theta_1 + \theta_2)^2} + \underline{\theta_3^2} = \|\theta\|_F^2$

$$K(x, x') = \exp(-\frac{1}{2k^2} \underline{\underline{\|x - x'\|^2}}).$$

$$K = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\theta^T K \theta = [\theta_1 \ \theta_2 \ \theta_3] \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$= (\theta_1 + \theta_2)^2 + \theta_3^2$$

$$\min_{\theta} \|Y - K\theta\|_2^2 + \lambda \theta^T K \theta,$$

$$\Rightarrow \hat{\theta} = (K + \lambda I)^{-1} Y.$$

Classification. $\{x^{(i)}, y^{(i)}\}$
 $y^{(i)} \in \{\pm 1\}$.

$$L(\theta) = \sum_{i=1}^n \sigma(y^{(i)} f_\theta(x^{(i)}))$$

$$= \sum_{i=1}^n \sigma(y^{(i)} \sum_{j=1}^n \theta_j K(x^{(i)}, x^{(j)}))$$

$$+ \lambda \theta^T K \theta$$

Kernel = infinite basis function

$$\{\psi_\ell(x) : \ell = 0, 1, 2, \dots, \infty\}$$

Example : $\psi_\ell(x) = \underline{x^\ell}$

$$f(x) = \sum_{\ell=0}^{\infty} w_\ell \underline{\psi_\ell(x)}$$

$$L(\underline{w}) = \hat{E}_D \left[(y - \sum_{\ell=0}^{\infty} w_\ell \underline{\psi_\ell(x)})^2 \right]$$

$$+ \lambda \cdot \sum_{\ell=0}^{\infty} \frac{w_\ell^2}{\psi_\ell(x)}$$

$\lambda > 0$ Regularization Coeff.

for w_ℓ .

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \underline{\lambda} \geq 0$$

$$\lambda_\ell \underset{\text{Small}}{\approx} w_\ell \rightarrow 0$$

Define

$$\hat{w} = \arg \min_w E_p \left[y - \sum_{\ell=0}^{\infty} w_\ell f_\ell(x) \right]^2 + \sum_{\ell=0}^{\infty} \frac{w_\ell^2}{\lambda_\ell}$$

$$\hat{f}(x) = \sum_{\ell=0}^{\infty} \hat{w}_\ell f_\ell(x)$$

$$\text{Then } \hat{f}(x) = \sum_{i=1}^n \theta_i k(x, x_i).$$

for θ_i , and

$$K(x, x') = \sum_{\ell=0}^{\infty} \lambda_\ell \underline{Y_\ell(x)} \underline{Y_\ell(x')}$$

And.

$$\sum_{\ell=0}^{\infty} w_\ell / \lambda_\ell = \underline{\theta^T K \theta}$$

$$\hat{\theta} = \operatorname{argmin}_{\theta} E_D \left[(y - \sum_{i=1}^n \theta_i k(x, x^{(i)}))^2 \right] + \underline{\theta^T K \theta}$$

Proof:

$$\min_w E_D \left[(y - \sum_{\ell=0}^{\infty} w_\ell Y_\ell(x))^2 \right] + \sum_{\ell=0}^{\infty} \frac{w_\ell^2}{\lambda_\ell}$$

$$\Rightarrow \text{Define } \underline{\delta(x)} = y - \sum_{\ell=0}^{\infty} \underline{\omega_\ell} \underline{\varphi_\ell(x)}$$

$$\nabla_{\omega_\ell} L(\hat{\omega}) = 2E_D[-\delta(x) \varphi_\ell(x)] + \frac{2\hat{\omega}_\ell}{\lambda_\ell} \geq 0$$

$$\Rightarrow \hat{\omega}_\ell = \lambda \hat{E}_D[\delta(x) \varphi_\ell(x)]$$

$$= \frac{\lambda}{n} \sum_{i=1}^n [\delta(x^{(i)}) \varphi_\ell(x^{(i)})]$$

$$f(x) = \sum_{\ell=0}^{\infty} \hat{\omega}_\ell \underline{\varphi_\ell(x)}$$

$$= \sum_{\ell=0}^{\infty} \frac{\lambda_\ell}{n} \sum_{i=1}^n \delta(x^{(i)}) \underline{\varphi_\ell(x^{(i)})} \underline{\varphi_\ell(x)}$$

$$= \sum_{\ell=0}^{\infty} \sum_{i=1}^n \frac{\lambda_\ell}{n} \delta(x^{(i)}) \underline{\varphi_\ell(x^{(i)})} \underline{\varphi_\ell(x)}$$

$$= \sum_{i=1}^n \sum_{l=0}^{\infty} \frac{\lambda_l}{n} \delta(x^{(i)}) \varphi_l(x^{(i)}) \varphi_l(x)$$

$$= \sum_{i=1}^n \frac{\delta(x^{(i)})}{n} \left[\sum_{l=0}^{\infty} \lambda_l \varphi_l(x^{(i)}) \varphi_l(x) \right]$$

$$= \sum_{i=1}^n \frac{\delta(x^{(i)})}{n} K(x^{(i)}, x)$$

$$= \left[\sum_{i=1}^n \Theta_i K(x^{(i)}, x) \right]$$

$$\boxed{\Theta_i \triangleq \frac{\delta(x^{(i)})}{n}}$$

Two Views :-

① $\hat{f}(x) = \sum_{i=1}^n \hat{\theta}_i k(x^{(i)}, x)$

② $\hat{f}(x) = \sum_{\ell=0}^{\infty} \hat{\omega}_{\ell} \varphi_{\ell}(x)$

\Leftrightarrow $k(x, x') = \sum_{\ell=0}^{\infty} \hat{\omega}_{\ell} \varphi_{\ell}(x) \varphi_{\ell}(x')$.

$\hat{\omega}_{\ell} \geq 0$

① If $k(x, x')$ is positive definite.

$\Leftrightarrow K = [k(x^{(i)}, x^{(j)})]_{i,j=1}^n \geq 0$.

for any n , $\{x^{(i)}\}$

$$\textcircled{2} \Leftrightarrow \exists \psi_\ell, \lambda_\ell \geq 0$$

Gaussian RBF Kernel:

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

$$= \exp(-\gamma \|x\|^2 + 2x^T x' - \gamma \|x'\|^2)$$

$$= \underbrace{\exp(-\gamma \|x\|^2)}_{\cdot} \underbrace{\exp(-\gamma \|x'\|^2)}_{\cdot} \underbrace{\exp(2x^T x')}_{\text{!!!}}$$

$$\exp(t) = 1 + t + \frac{t^2}{2} + \dots$$

$$= \sum_{l=0}^{\infty} \frac{t^l}{l!}$$

$$\exp(2x^T x') = \boxed{\sum_{l=0}^{\infty} \frac{(2x^T x')^l}{l!}}$$

⇒ Assume zero.

$$K(x, x') = \sum_{l=0}^{\infty} \lambda_l \varphi_l(x) \varphi_l(x')$$

$$\varphi_l(x) = \exp(-\gamma x^2) x^l.$$

$$\lambda_l = \frac{(2\gamma)^l}{l!} \geq 0$$

