

HW 2 : Regression using kernels

Presentation paper] ^{Last} 2 Wed of "teaching" period
Project (presentation + 1-page proposal + report (5 pages))

Regression

Linear regression using linear model

$$h_{\omega}(x) = x \cdot \omega$$

$\{(x_i, y_i)\}_{i \in [m]}$ iid from \mathcal{D}

$$\text{Solve ERM: } \min_{\omega \in \mathbb{R}^d} \frac{1}{m} \sum_{i \in [m]} \|y_i - h_{\omega}(x_i)\|^2$$

$$= \min_{\omega \in \mathbb{R}^d} \|X\omega - Y\|^2$$

$$X \in \mathbb{R}^{m \times d} \quad Y \in \mathbb{R}^m$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \quad (\nabla_{\omega} \|X\omega - Y\|^2 = 0)$$

Solution using optimization:

$$\omega = (X^T X)^{-1} X^T Y$$

(assuming that $X^T X$ is invertible)

Interpretation of LS regression from linear algebraic viewpoint

- when X is square, $m = d$ and X is invertible,

$$\omega = X^{-1} Y$$

$$X^{-1} Y = (X^T X)^{-1} X^T Y$$

- $m > d$ overdetermined case: $X\omega = Y$ may not have a solution!

$$\begin{bmatrix} X \\ Y \end{bmatrix} \begin{bmatrix} \omega \end{bmatrix} = \begin{bmatrix} Y \end{bmatrix} \quad \text{But, } X\omega = Y \text{ has a solution when } Y \text{ is in the } \text{Ran}(X)$$

$$X^T X \omega = X^T Y \Rightarrow \omega = (X^T X)^{-1} X^T Y$$

if $X^T X$ is invertible, that is X has full rank = d ,

- underdetermined / overparameterized $d > m$

$$\begin{bmatrix} X \\ Y \end{bmatrix} \begin{bmatrix} \omega \end{bmatrix} = \begin{bmatrix} Y \end{bmatrix}$$

when does $X\omega = Y$ have a solution?

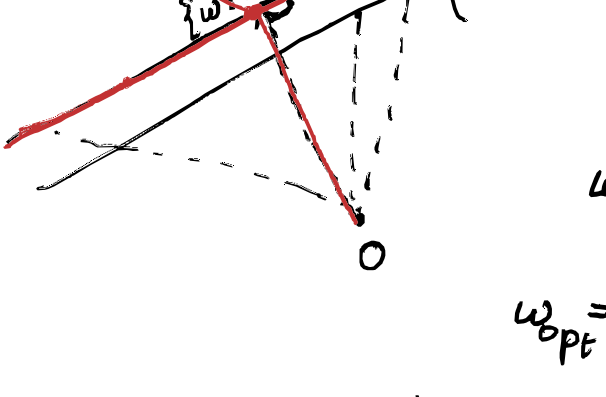
when X^T has full rank = m .

Is solution unique? No.

if $z \in \text{Nul}(X)$ and ω solves $X\omega = Y$, then $X(\omega + z) = Y$

Mostly, we

Want "minimum norm" solution for ω .



$$X\omega = Y$$

$$\omega_{\text{opt}} \in \text{Ran}(X^T)$$

$$\omega_{\text{opt}} = X^T z$$

$$X\omega = X X^T z = Y$$

$$z = (X X^T)^{-1} Y$$

($X X^T$ is invertible)

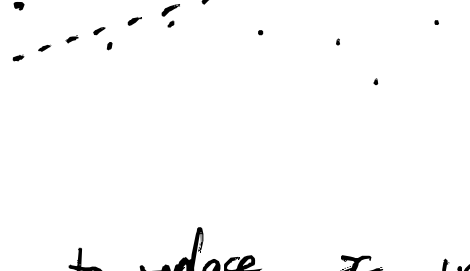
when X has full row rank = m , $X X^T$ is positive def

$$\omega_{\text{opt}} = X^T z$$

$$= X^T (X X^T)^{-1} Y$$

if no ω exists st. $X\omega = Y$

$$\arg \min_{\omega} \|X\omega - Y\|^2 = (X^T X)^{-1} X^T Y$$



Want: to replace x with some "features" $\Phi(x)$

$$X = \begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_m) \end{bmatrix}$$

Solution of LS regression:

$$\omega = (X^T X)^{-1} X^T Y$$

o/p fn: $h_{\omega}(x) = \omega \cdot \Phi(x)$

Questions: • what are features $\Phi(x)$?

• why/when is lifting linear regression to feature space effective?

$h_{\omega}(x)$: still a linear function on feature space (of $\Phi(x)$)

but can be complicated nonlinear fn on data space (\mathcal{X})

Reproducing kernel Hilbert space

RKHS

$\mathcal{X} \subseteq \mathbb{R}^d$ compact

Theorem: If $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is PD,

$\exists!$ RKHS(\mathcal{H}) which is defined by the following properties:

$$\| \cdot \|^2 = \langle \cdot, \cdot \rangle$$

- \exists some inner product s.t.

$$\langle \kappa(x, \cdot), \kappa(x', \cdot) \rangle = \kappa(x, x')$$

$$\mathcal{H} = \overline{\text{span} \{ \kappa(x, \cdot) : x \in \mathcal{X} \}}$$

$$\text{for any } f \in \mathcal{H}, \quad \langle f, \kappa(x, \cdot) \rangle = f(x)$$

Recall: κ is PD kernel if Gram matrix

$$G_{ij} = \kappa(x_i, x_j) \quad i, j \in N \text{ is SPSD}$$

Reproducing property

Evaluation functional

$$E_x: \mathcal{H} \rightarrow \mathbb{R}$$

$$E_x \in \mathcal{H}^*$$

$$E_x f = f(x)$$

$$E_x(f + g) = f(x) + g(x) = E_x f + E_x g$$

Linearity \checkmark

$$|E_x f| = |f(x)| \stackrel{\text{suppose}}{<} \infty \quad \left(\begin{array}{l} \text{e.g. } f \text{ is} \\ \text{continuous,} \\ \mathcal{X} \text{ is compact} \end{array} \right)$$

if E_x (evaluation functional) at every x is bounded and linear, then,

(or $|E_x f| < C \|f\|_{\mathcal{H}}$), Riesz representation theorem holds

So, $\exists!$ $g_x \in \mathcal{H}$ s.t.

$$\langle f, g_x \rangle_{\mathcal{H}} = E_x f = f(x).$$

$$g_x \equiv \kappa(x, \cdot)$$

\mathcal{X}

Proof:

Want: there is an RKHS corresponding to a PD κ .

$$\mathcal{H}_0 = \text{span}\{x(x_i, \cdot)\}_{i \in \mathcal{I}} \quad \mathcal{I} \text{ is finite}$$

$$f \in \mathcal{H}_0 \Rightarrow f(x) = \sum_{i \in \mathcal{I}} \alpha_i x(x_i, x)$$

Define

$$\langle f, g \rangle := \sum_{i, j \in \mathcal{I}} \alpha_i \beta_j x(x_i, x_j)$$

$$\text{if } f = \sum_i \alpha_i x(x_i, \cdot) \quad g = \sum_j \beta_j x(x_j, \cdot)$$

Want: show \langle, \rangle is an inner product.

- bilinearity:
 - symmetric $\langle f, g \rangle = \langle g, f \rangle$

$$\langle f + g, h \rangle = \langle f, h \rangle + \langle g, h \rangle$$

$$\alpha \langle f, h \rangle = \langle \alpha f, h \rangle$$

- positive definite

$$\langle f, f \rangle \geq 0$$

$$\text{and } \langle f, f \rangle = 0 \text{ iff } f = 0$$

$$\langle f, f \rangle \geq 0$$

$$\Rightarrow \sum_{i, j \in \mathcal{I}} \alpha_i \alpha_j x(x_i, x_j) = \alpha^T \underbrace{G}_{\text{Gram matrix}} \alpha$$

$$\geq 0$$

$$\text{because } G \text{ is SPSPD. } (\because \kappa \text{ is PD})$$

- if $f = 0$, then, $\langle f, f \rangle = 0$ ($\because \alpha_i = 0$).

- if $\langle f, f \rangle = 0$, then $f = 0 \in \mathcal{H}_0$.

$$\langle f, f \rangle = \sum \alpha_i \alpha_j x(x_i, x_j)$$

Need to use Cauchy Schwartz for PD kernels!

Lemma: if κ is a PD kernel,

$$\text{then, } |\kappa(x, x')|^2 \leq |\kappa(x, x) \kappa(x', x')|$$

Proof:

$$G = \begin{bmatrix} \kappa(x, x) & \kappa(x, x') \\ \kappa(x', x) & \kappa(x', x') \end{bmatrix} \quad \begin{aligned} & \kappa(x, x') \cdot \kappa(x', x) \\ & = (\kappa(x, x'))^2 \\ & (\because \kappa \text{ is sym}) \end{aligned}$$

$$\kappa(x, x) \kappa(x', x') - \kappa(x, x')^2 \geq 0$$

$$\kappa(x, x') = \langle \phi(x), \phi(x') \rangle$$

$$|\langle \phi(x), \phi(x') \rangle|^2 \leq \|\phi(x)\|^2 \|\phi(x')\|^2$$

- Define $\tau(f, f) = \langle f, f \rangle$ is PD kernel on \mathcal{H}_0 .

$$G_{ij} = \tau(f_i, f_j) \quad f_i, f_j \in \mathcal{H}_0$$

$$\sum_{i, j} c_i c_j G_{ij} = \sum_{i, j} c_i c_j \tau(f_i, f_j)$$

$$= \sum_{i, j} c_i c_j \langle f_i, f_j \rangle$$

$$= \langle \sum_i c_i f_i, \sum_j c_j f_j \rangle$$

$$(\text{linearity of inner product}) \geq 0$$

$$\text{So, } G \text{ is SPSPD.}$$

- CS for τ .

$$|\tau(f, x(x, \cdot))| \leq \tau(f, f) \tau(x(x, \cdot), x(x, \cdot))$$

$$|\langle f, x(x, \cdot) \rangle| \leq \langle f, f \rangle x(x, x)$$

Use reproducing property.

$$|f(x)| = |\langle f, x(x, \cdot) \rangle|$$

Use CS for τ

$$|f(x)| = |\langle f, x(x, \cdot) \rangle|$$

$$\leq \langle f, f \rangle \underbrace{\langle x(x, \cdot), x(x, \cdot) \rangle}_{\text{repr. prop.}}$$

$$= \langle f, f \rangle x(x, x)$$

$$\therefore \text{ if } \langle f, f \rangle = 0, \text{ then, } |f(x)| = 0$$

$$\text{at all } x$$

$$\Rightarrow f = 0 \in \mathcal{H}_0$$

We have shown that

$$\langle f, g \rangle := \sum_{i, j} \alpha_i \beta_j x(x_i, x_j)$$

$$\text{where } f = \sum \alpha_i x(x_i, \cdot)$$

$$g = \sum \beta_j x(x_j, \cdot)$$

is a valid inner product.

$$\Rightarrow \mathcal{H}_0 \text{ is a pre-Hilbert space}$$

$$\mathcal{H}(\mathcal{H}_0) \overset{\mathcal{H}}{=} \text{is a Hilbert space}$$

$$\mathcal{H}_0 \cup \{\text{limit points}\}$$

$$\mathcal{H}_0 \text{ is dense in } \mathcal{H}$$

by Hahn-Banach Theorem, reproducing prop. holds in \mathcal{H} .

Heine-Borel: closed + bounded sets of $\mathbb{R}^n \equiv \text{compact}$

If there are 2 kernels κ, τ associated with \mathcal{H} ,

$$\langle \kappa(x, \cdot), \tau(y, \cdot) \rangle = \tau(y, x)$$

$$= \kappa(y, x)$$

$$\Rightarrow \kappa \text{ is unique.}$$

Next time: Mercer map.

Reproducing property

$$\langle f, g \rangle = \sum_{i,j} \alpha_i \beta_j \kappa(x_i, x_j)$$

if $g = \kappa(x_j, \cdot)$ $f(x) = \sum_{i \in I} \alpha_i \kappa(x_i, x)$

then,

$$\begin{aligned} \langle f, \kappa(x_j, \cdot) \rangle &= \sum_i \alpha_i \kappa(x_i, x_j) \\ &= f(x_j) \end{aligned}$$

satisfied!