

Recap: With (pro. our $S \sim \mathcal{D}^m$), $\geq 1-\delta$
 $R(h_S) \leq \frac{1}{m} \log \frac{|\mathcal{H}|}{\delta}$ (\mathcal{D}^m)
 \uparrow
ERM over training S $m = |S|$
 $|\mathcal{H}|$: size of finite hypothesis class \mathcal{H} .

$$R(h_S) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(x,y, h_S)$$

More on generalization bounds after we study regression & kernel regression

Linear regression

$$h(x) = w^T x + b$$

$$\ell(x, y, h) = (w^T x + b - y)^2$$

$$\hat{R}_S(h) = \sum_{(x,y) \in S} \ell(x, y, h)$$

MSE

Solving ERM over linear models (w, b) is called linear regression

Want w, b s.t.

$\hat{R}_S(h)$ is mini
 \hookrightarrow parameterized by w, b .

$$w \equiv [w, b] \in \mathbb{R}^{d+1}$$

$$x \equiv [x, 1] \in \mathbb{R}^{d+1}$$

$$\min_w \hat{R}_S(h) = \min_w \frac{1}{m} \sum_{(x,y) \in S} \|y - w \cdot x\|^2$$

WKT at minimum w^* ,

$$\nabla_w \hat{R}_S(h) = 0$$

$$\Rightarrow \sum_{(x,y) \in S} (y - w \cdot x) x = 0 \quad (*)$$

Define $Y = [y_1, \dots, y_m]^T \in \mathbb{R}^m$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}$$

$(*)$ if $X^T X$ is invertible
 $X^T Y = X^T X w \Rightarrow w = (X^T X)^{-1} X^T Y$

$$X^T (Y - Xw) = 0 \in \mathbb{R}^{d+1}$$

Linear algebra

$$\min_w \hat{R}_S(w) = \frac{1}{m} \sum_{i=1}^m (y_i - w^T x_i)^2$$

$$= \min_w \frac{1}{m} \|Y - Xw\|^2$$

Least squares problem

Case 1: $Y = Xw \Rightarrow w = X^{-1} Y$
 $d+1 = m$ and X is invertible

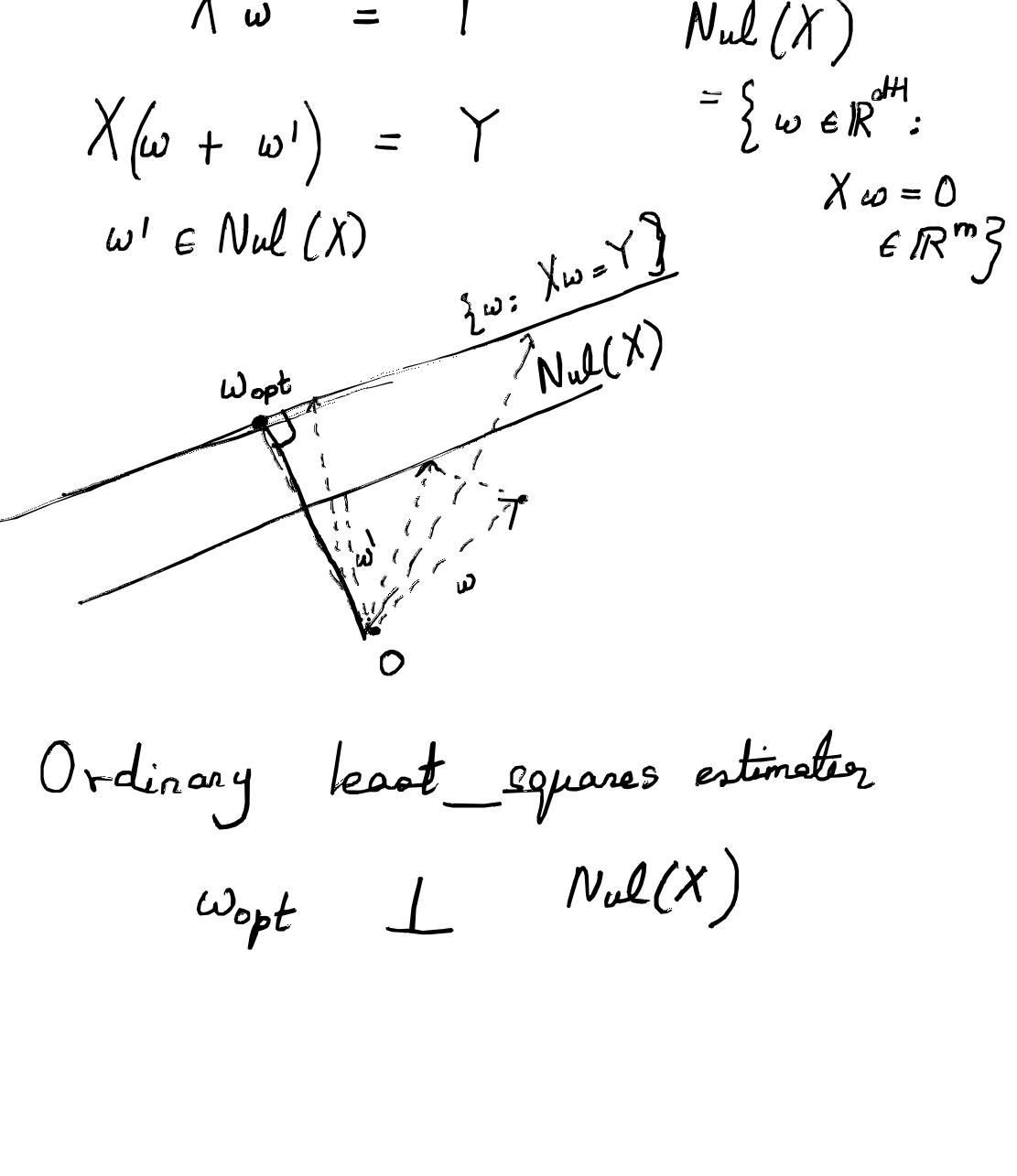
Case 2: $d+1 > m$ overparameterized non-unique!
 $XX^T \in \mathbb{R}^{m \times m}$ solution for $Xw = Y$.
if X has full row rank ($\text{rank} = m$)

$$Y = Xw$$

$$\boxed{w = X^T (XX^T)^{-1} Y} \text{ is optimal in norm}$$

$$Xw = XX^T (XX^T)^{-1} Y = Y$$

ie min-norm solution. why?



Ordinary least squares estimator

$$w_{\text{opt}} \perp \text{Nul}(X)$$

$$X w_{\text{opt}} = Y$$

$$\text{Nul}(X) \perp \text{Ran}(X^T)$$

Min-norm solution

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

$$X w_{\text{opt}} = XX^T z = Y$$

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

$$w_{\text{opt}} = X^T (XX^T)^{-1} Y$$

$$w_{\text{opt}} \perp \text{Nul}(X)$$

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

when X is full row rank ($\text{rank} = m$)
 $m < d+1$

Case 3: $m > d+1$ our determined rank ($d+1$)
 X has full column rank
then $X^T X$ is invertible and

$$w = (X^T X)^{-1} X^T Y \text{ (Normal equation solution)}$$

also solves $Xw = Y$ if $Y \in \text{Ran}(X)$

More generally, $w = (X^T X)^{-1} X^T Y$
is only OLS solution
(see optimization approach above)

Soln to linear regression

$$h(x) = \omega^T x \quad \text{linear model}$$

$$x \rightarrow \Phi(x)$$

$$h(x) = \omega^T \Phi(x) + b$$

$$= \langle \omega, \Phi(x) \rangle \rightarrow \text{inner product in function space}$$

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

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}$$

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

$$X^T X_{ij} = \langle \Phi(x_i), \Phi(x_j) \rangle$$

$$m \times m \quad = \quad \chi(x_i, x_j) \quad \text{for some PD kernel}$$

$$X^T Y_{m \times 1} = \langle \Phi(x_i), \Phi(y_i) \rangle$$

To compute ω_{opt} , we only need to invert $(X^T X)$: Gram matrix

Linear regression \rightarrow kernel regression

if $X^T X$ is not invertible

$$\omega_{\text{opt}} = (X^T X + \underset{\substack{\uparrow \\ \text{regularization} \\ \text{parameter}}}{\lambda} I)^{-1} X^T Y$$

$$x \in \mathcal{X} \quad \Phi(x) \in \mathbb{R}^D$$

" or \mathcal{H} .

d-dim.

Many feature maps can be associated with the same kernel!

Reproducing kernel Hilbert space RKHS

For any PD kernel $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$,
 \exists a corresponding Hilbert space \mathcal{H}
 s.t. for any $f \in \mathcal{H}$

$$f(x) = \langle f, \kappa(\cdot, x) \rangle_{\mathcal{H}}$$

Reproducing property

Consider some arbitrary data points
 x_1, x_2, \dots, x_m

$\text{span} \{ \kappa(\cdot, x_i) \}$

$f, g \in \text{span} \{ \kappa(\cdot, x_i) \}$

$$f(x) = \sum_{i=1}^m \alpha_i \kappa(x, x_i) \leftarrow$$

$$g(x) = \sum_{i=1}^m \beta_i \kappa(x, x_i)$$

Define

$$\langle f, g \rangle := \sum_{i,j=1}^m \alpha_i \beta_j \kappa(x_i, x_j)$$

Linearity: Bilinear

$$\begin{aligned} \langle f, g \rangle &= \sum_{j=1}^m f(x_j) \beta_j \\ &= \sum_{i=1}^m \alpha_i g(x_i) \end{aligned}$$

$$\begin{aligned} \langle f_1 + f_2, g \rangle &= \sum_{i=1}^m (\alpha_i^{(1)} + \alpha_i^{(2)}) g(x_i) \\ &= \sum_{i=1}^m \alpha_i^{(1)} g(x_i) + \sum_{i=1}^m \alpha_i^{(2)} g(x_i) \\ &= \langle f_1, g \rangle + \langle f_2, g \rangle \end{aligned}$$

Positive definiteness

$$\begin{aligned} \langle f, f \rangle &= \sum_{i,j=1}^m \alpha_i \alpha_j \kappa(x_i, x_j) \\ &= \alpha^T G \alpha \geq 0 \end{aligned}$$

(G is Gram matrix, which is SPSD for PD kernel κ)

if $f = 0$, $\langle f, f \rangle = 0$

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

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

To show: completion under $\langle \cdot, \cdot \rangle$

$$\begin{aligned} \langle \kappa(\cdot, x_i), \kappa(\cdot, x_j) \rangle_{\mathcal{H}} \\ = \kappa(x_i, x_j) \end{aligned}$$

$$\rightarrow | \langle \kappa(\cdot, x_i), \kappa(\cdot, x_j) \rangle_{\mathcal{H}} |^2 \leq \|\kappa(\cdot, x_i)\|_{\mathcal{H}}^2 \|\kappa(\cdot, x_j)\|_{\mathcal{H}}^2$$

PD kernel

$$\|\kappa(\cdot, x_i)\|_{\mathcal{H}}^2 := \langle \kappa(\cdot, x_i), \kappa(\cdot, x_i) \rangle$$

$\langle \cdot, \cdot \rangle_{\mathcal{H}}$ inner product satisfies CS.

Next: RKHS feature map