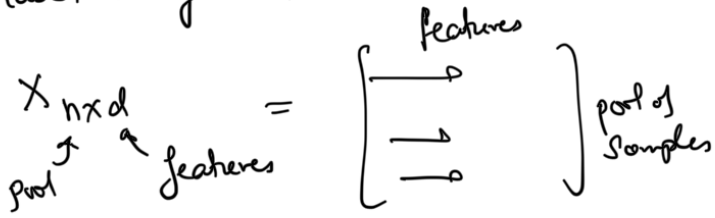


Experimental Design

find k points out of n that $\min_{|S| \leq k} R(\hat{f}) - R(f^*)$ es. $f^* = \theta^{*T} x$
 $\{x_1, \dots, x_k\} = S$ $Y = \theta^{*T} x + \epsilon$

label budget k



Regression $R(\hat{f}) - R(f^*) = E_{x, y, D} [(y - \hat{f}(x))^2] - E_{x, y} [(y - f^*(x))^2]$
← data (training) σ^2 Irreducible error

$$E[(y - \hat{f}(x))^2] = E_{x, y} [(y - f^*(x))^2] + E_{x, D} [(f^*(x) - \hat{f}(x))^2]$$

$y = f^*(x) + \epsilon$
 $\epsilon \sim \text{iid } N(0, \sigma^2)$

$$= E[(y - f^*(x))^2] + E_x [(f^*(x) - E_D[\hat{f}(x)])^2] + E_{x, D} [(\hat{f}(x) - E_D[\hat{f}(x)])^2]$$

Cross-terms vanished 1 $E[y] f^*(x) = E[\hat{f}(x)] = E[f^*(x)]$

$$= R(f^*) + \underbrace{\text{bias}^2}_{\sigma^2 \text{ irreducible}} + \text{variance}$$

$$(2 E[(f^*(x) - E_D[\hat{f}(x)]) E_D[\hat{f}(x) - E_D[\hat{f}(x)]]])$$

Bias can't be improved by sampling



\Rightarrow choice of sampling can only improve variance. ← measure of uncertainty

\therefore Exp design (1-shot) & active learning (sequential)
 minimize variance.

Goal: How to choose x_1, \dots, x_k to minimize $E_{x, D} [(\hat{f}(x) - E_D[\hat{f}(x)])^2]$?

$$= E[\langle \hat{\theta} - \theta^*, x \rangle^2]$$

$\hat{f}(x) = \hat{\theta}^T x$ $\hat{\theta} = (X^T X)^{-1} X^T Y$ LS

$$y = \theta^T x + \varepsilon \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$E[\hat{\theta}] = \theta^*$$

Uncertainty / variance of prediction

$$\rightarrow E_{\mathcal{D}}[\langle \hat{\theta} - \theta^*, x \rangle^2] = x^T \underbrace{E_{\mathcal{D}}[(\hat{\theta} - \theta^*)(\hat{\theta} - \theta^*)^T]}_{\text{Cov}(\hat{\theta})} x \quad \leftarrow$$

$$\begin{aligned} \text{Cov}(\hat{\theta}) &= E[(X^T X)^{-1} X^T \varepsilon \varepsilon^T X (X^T X)^{-1}] \\ &= (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \quad \leftarrow \end{aligned}$$

$$\theta^* = (X^T X)^{-1} X^T y$$

$$E[\varepsilon \varepsilon^T] = \sigma^2 I$$

A-optimality

$$\min_{|S| \leq k} \text{var}(\hat{\theta}_S)$$

$$= \sigma^2 \text{tr}[(X_S^T X_S)^{-1}]$$

$S \subseteq \mathcal{D}$

E/V-optimality

$$\min_{|S| \leq k} \text{var}(\langle \hat{\theta}_S, x \rangle)$$

$$= \sigma^2 x^T (X_S^T X_S)^{-1} x \quad \leftarrow$$

Transfer

Combinatorial opt / Intractable

$$\min_{|S| \leq k} f(X_S^T X_S) \equiv f(\Sigma)$$

	A-opt	$\text{tr}(\Sigma^{-1})$	Average	
	E-opt	$\ \Sigma^{-1}\ _2$	Eigenvalue	$\max_{\ x\ \leq 1} x^T M x = \lambda_{\max}$
max volume of region covered by chosen points \leftarrow	D-opt	$\det(\Sigma^{-1})$	Determinant	\leftarrow tractable $-\log(\ \Sigma\)$
	T-opt	$(\text{tr}(\Sigma))^{-1}$	Trace	Complex
	V-opt	$\text{tr}(X \Sigma^{-1} X^T)$	Variance (In-sample)	
	G-opt	$\max_{\text{diag}}(X \Sigma^{-1} X^T)$	Global	

$$\underline{\underline{X_S^T X_S}} = X^T W X \quad W = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & & \\ \vdots & & \ddots & \\ 0 & & & 1 \end{bmatrix} = \text{diag}(w) \quad w_i \in \{0,1\}$$

$\sum_i w_i \leq k$

Continuous Relaxation

$$\begin{aligned} & X^T W X \\ & \uparrow \\ \rightarrow & \min_{\substack{w_i \in [0,1] \\ \sum_i w_i \leq k}} \text{tr}(X^T W X) = \sum_i w_i x_i x_i^T \\ & \frac{w_i}{k} \equiv p_i \quad \sum_i \frac{w_i}{k} = \sum_i p_i \leq 1 \end{aligned}$$

Thm: $f\left(\sum_{i \in S} w_i x_i x_i^T\right) \leq \min_{|S| \leq k} f(X_S^T X_S)$

\geq

How to choose subset X_S given $w_i \in [0,1]$?

Can we still get $f(X_S^T X_S) \leq (1+\epsilon) \min_{|S| \leq k} f(X_S^T X_S)$?

→ Treat $\{\frac{w_i}{k}\}$ as p_i + sample with or without replacement k points

Thm: Approximation guarantee

rough A, V opt

$(1+o(1))$

$k = \Omega(d)$

$1+\epsilon$

$k = \Omega\left(\frac{d \log d}{\epsilon^2}\right)$ or $k = \Omega\left(\frac{d^2}{\epsilon}\right)$ without replacement

More sophisticated rounding methods can get $1+\epsilon$ approximation for $k = \Omega\left(\frac{d}{\epsilon}\right)$ for all optimality criteria.

Show for G-optimality $O(d^2)$ points suffice without replacement to get approximation ~~to~~ ratio of 1.


Fact: For every w s.t. $\sum_i w_i = 1, w_i \geq 0$

$$\exists w' \text{ s.t. } \sum_i w'_i = 1, w'_i \geq 0 \text{ with } \|w'\|_0 = O(d^2).$$

$$\text{s.t. } \underline{X^T W X} = X^T W' X$$

Caratheodory's Thm: If a point z lies in convex hull of set $\Omega \in \mathbb{R}^m$ then z can be written as convex combⁿ of at most $m+1$ extremal points.

Let $\Omega = \{ \text{vec}(x_i x_i^T) \}_{i=1}^n \in \mathbb{R}^{d^2}$



$z = \text{vec}(\sum_i w_i x_i x_i^T)$ lies in convex hull of Ω .

Gre-optimal procedure (with replacement)

Input: $X_{n \times d}$, budget k

1. Solve relaxed problem to get w .
2. Has $O(d^2)$ sparse solution w' with same objective as w .
3. Observe data pt. i $n_i = \lceil w_i k \rceil$ times (collect Y_i)
4. Compute $\hat{\theta} = (X^T X)^{\dagger} X^T Y$

Thm: w.p. $\geq 1-\delta \quad \forall x \in \mathcal{X}$

$$\begin{aligned} |x^T \hat{\theta} - \theta^*| &\leq \sqrt{\|x^T (X^T X)^{\dagger} x\|} \sqrt{2\sigma^2 \log(2|\mathcal{X}|/\delta)} \\ &= \sqrt{\frac{2d\sigma^2 \log(2|\mathcal{X}|/\delta)}{k}} \end{aligned}$$

Proof: $x^T(\hat{\theta} - \theta^*) = x^T (X^T X)^{\dagger} X^T \varepsilon =: z^T \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

$$\sim \mathcal{N}(0, \sigma^2 x^T (X^T X)^{\dagger} x)$$

$$\Pr(|x^T(\hat{\theta} - \theta^*)| > \varepsilon) \leq 2e^{-\frac{\varepsilon^2}{2\sigma^2 x^T (X^T X)^{\dagger} x}}$$

Union bound $\forall x \in \mathcal{X}$

$$\Pr(|\varepsilon| > 0) \leq 2e^{-\frac{\varepsilon^2}{2\sigma^2}}$$

$$\omega_p \geq 1 - \delta$$

$$P(A \cup B) \leq P(A) + P(B)$$

$$|x^T(\hat{\theta} - \theta^*)| \leq \sqrt{2\sigma^2 \log \frac{2|x|}{\delta}} \cdot \underbrace{\sqrt{x^T (X^T X)^{-1} x}}_{\sqrt{\frac{d}{k}}} \quad \curvearrowright$$