Recap Exp design

Choosing $\underbrace{x_1 \ldots x_k}_{S}$ out of $n$ s.t. $\underset{=}{R(\hat{f}_s) - R(f^*)}$ ... $y \approx f^*(x) + \varepsilon$, $\theta^{*T} x$, $\varepsilon \sim$ iid $N(0, \sigma^2)$

linear regression

$E[<\hat{\theta}_s - \theta^*, x>]$

$\underline{x^T (X_s^T X_s)^{-1} x} \Leftarrow$ predictor var

A-opt $\quad E[<\hat{\theta}_s - \theta>^2] \equiv (X_s^T X_s)^{-1} \quad$ var of parameters

E-opt $\quad \underset{x}{\max}\, x^T (X_s^T X_s)^{-1} x$

V-opt $\quad X^T (X_s^T X_s)^{-1} X$

$$f((X_s^T X_s)^{-1}) = f((X^T W X)^{-1})$$

combinatorial opt $\quad W = \begin{bmatrix} 0 & 1 & & 0 \\ & 0 & 1 & \\ 0 & & & \ddots \end{bmatrix} = diag(w)$, $w_i \in \{0,1\}$

$\sum w_i = k$

Continuous relaxation

$\quad w_i \in [0,1] \quad \Leftarrow$

$\exists$ a ~~scaled~~ vector $w'$ that is $O(d^2)$ sparse s.t. $x^T W x = x^T W' x$

$\sum w_i' = 1 \quad \sum w_i = 1$

Optimal exp design procedure G-opt
_____

1. Solve cont relaxed version of G-opt to get $w$.

2. Has $O(d^2)$ sparse $w'$

3. Sample $i^{th}$ data points $n_i = \lceil w_i' k \rceil$ times $(x_i, y_i)$

4. Build est $\hat{\theta}$ using these pts.

Thm: $wp \geq 1 - \delta \quad \forall\, x \in \mathcal{X} \quad |<x, \hat{\theta} - \theta^*>| \leq \underbrace{\sqrt{||x (X_s^T X_s)^{-1} x||}}_{O(d/k)} \sqrt{2\sigma^2 \log 2|X|/\delta}$

$\simeq \sqrt{\frac{d}{k}}$

where $\hat{\theta}$ constructed using

... $1 + w_i' k = O(d^2) + k \Leftarrow$

$$\sum_i n_i = \sum_{i \sim O(d^2)} \lfloor w_i k \rfloor \leq \underset{i \sim O(d^2)}{2} \quad \quad [\because \sum w_i' = 1]$$

data points/labels.

$$\bigstar \quad x^T (X_S^T X_S)^{\dagger} x = x^T \left( \sum_i n_i x_i x_i^T \right)^{\dagger} x \leq x^T \left( \sum \lfloor w_i' k \rfloor x_i x_i^T \right)^{\dagger} x$$
$$\underset{\sim O(d^2)}{}$$

$$\leq \frac{x^T \left( \sum_i w_i' k \, x_i x_i^T \right)^{\dagger} x}{k} \quad \leftarrow O(d)$$

$$\max_x x^T \left( \sum_i w_i' x_i x_i^T \right)^{\dagger} x \geq \sum_i w_i' x_i^T \left( X^T W X \right)^{\dagger} x_i = tr \left( X^T W X \left( X^T W X \right)^{\dagger} \right)$$

$$\text{max} \geq \text{avg} \quad \quad = tr(I) = d$$

Since $w'$ is optimal, achieves lower bound with equality $\supset$
$\quad\quad\quad\quad\quad$ (Kiefer-Wolfowitz Equivalence Theorem)

---

Better rounding techniques

$\quad 1+\varepsilon \quad$ using $k = \Omega\left(\frac{d^2}{\varepsilon}\right)$ or $\Omega\left(\frac{d}{a^2}\right) \quad \forall \, A,D,T,E,V,G$ $\quad\swarrow\overset{\text{regret}}{\underset{=}{\text{min}}}$

$\quad 1+\varepsilon \quad$ using $k = \Omega\left(\frac{d}{\varepsilon}\right) \quad$ only A,D opt. $\leftarrow$ Federov's also.
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ Greedy algo.

$\quad\quad\quad$ Federov - randomly sampling $k$ pts. &
$\quad\quad\quad\quad\quad\quad\quad$ swapping pts.

$\quad\quad\quad$ Greedy - small randomly samples.
$\quad\quad\quad\quad\quad\quad\quad$ greedy addition.

---

Nonlinear regression $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad{}^{y = f^*(x) + \varepsilon}\nearrow$

$\quad$ · Generalized linear regression $\quad g(f^*(x)) = \theta^{*T} x \quad\quad g\text{-form}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ known
$$f^*(x) = \frac{e^{\theta_0^T x}}{1 + e^{\theta_0^T x}}$$

$\quad\quad$ val of prediction $\quad \hat{\theta} - MLE$

$\quad\quad\quad$ under mild reg. $\quad MLE \quad E\left[\|\hat{\theta} - \theta^*\|^2\right] = (1 + o(1)) \, tr\left(\mathcal{I}(X, \theta^*)^{\dagger}\right)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\underline{\quad\quad}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\checkmark$

Fisher Information matrix

linear $I(X, \theta^*) = X^T X$

$$\min_{|S| \leq k} tr(I(X_S, \theta^*))^T$$

1. $k/2$ samples randomly $\to$ get estimate $\hat{\theta}_{MLE}$   $k/2$

2. plug-in

$$\tilde{X}_i \leftarrow func(X_i, \hat{\theta}, g)$$

- Neural Networks (deep)          Core-set sampling.

$$\left| \frac{1}{n} \sum_{i=1}^{n} \ell(x_i, y_i) - \frac{1}{|S|} \sum_{i \in S} \ell(x_i, y_i) \right| \leq Small$$

↑                              ↑

loss $\ell$ is Lipschitz            (don't want to use $\{y_i\}_{i=1}^n$)

$\equiv$ K-center problem

$$\min_{|S'| \leq k} \max_i \min_{j \in S' \cup S^o} \hat{d}(x_i, x_j)$$

↳    ↑     ↑        NP-hard

$S^o$ - initial dataset
$S'$ - new labeled dataset

greedy sol$^n$ - approx$^n$ ratio 2

---

Active Learning    - sequentially choose $x_1 \cdots x_T$ to min $R(\hat{f}_{X_1 X_T}) - R(f^*)$   $= *$

Regression  $x_t = \arg\min_x \hat{\sigma}_{t-1}(x)$

$\quad = $ closed form  linear models
                    generalized linear models
                    Bayesian models (GP)
                    ensembles for general nonlinear models
                                              e.g. NNs.

$*$  $\min_{x_t} R(\hat{f}_{x_1 \cdots x_t}) - R(f^*)$      $\hat{f}_1 \cdots \hat{f}_m$

$\quad y_1 \cdots y_{t-1}$                                     ↓

                                                    $x \to \hat{y}_1 \quad \hat{y}_m$

$$x \to \frac{1}{m}\sum_{j=1}^{m}\left(y_j - \frac{1}{m}\sum_{j=1}^{m}y_i\right)^2$$

## Classification

uncertainty of predicted labels

Binary classes $\quad Ber(p) \sim p(1-p) \quad$ max $p = \frac{1}{2}$

$$x \to P(Y=1|x), P(Y=0|x)$$

Logistic regression $\quad P(Y=1|x) = \dfrac{e^{\theta^T x}}{1+e^{\theta^T x}}$

$$P(Y=1|x) - \frac{1}{2} \quad \leftarrow \text{ near decision boundary}$$

Multiple classes

least confident ✓

$$\underset{x}{\arg\min} \quad 1 - \max_{y} P(y|x)$$
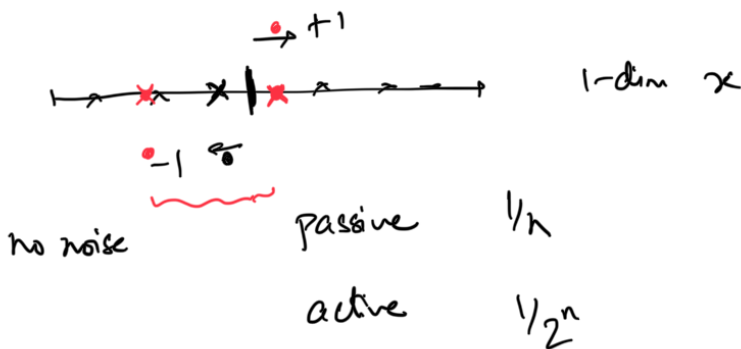
margin sampling ✓

$$\underset{x}{\arg\min} \; P(Y_{(1)}|x) - P(Y_{(2)}|x)$$

entropy sampling ✓

$$\underset{x}{\arg\max} \sum_{y} P(Y=y|x) \log \frac{1}{P(Y=y|x)}$$

How to extend these ideas to non-probabilistic classifiers?

Linear, SVM, Decision Trees

1-dim $x$

no noise $\qquad$ passive $\quad 1/n$

active $\quad 1/2^n$

d-dim    linear $\begin{cases} \text{passive} & d/n \\ \text{active} & e^{-n/d} \end{cases}$

no noise

with noise $\begin{cases} \text{passive} & \left(\frac{d}{n}\right)^{\frac{K}{2K-1}} \\ \text{active} & \left(\frac{d}{n}\right)^{\frac{K}{2K-2}} \end{cases}$

$\longrightarrow$

Algo    adaptive noise
        linear case.



Tsybakov
noise label