

Week 13

↳ Pset 12

- due Sunday (hard deadline)
- self-grade due Monday (hard deadline)

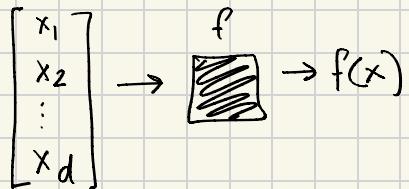
↳ Anyone interested in grading next semester?

↳ Midterm 11/25

Suggestion:

1. Review psets (including this week!)
2. Review recurring concepts (e.g., SVD, variance, etc)
3. Take 1st practice exam (75 min on paper). Review.
4. Take 2nd.

Explainable AI



Q: why did the model output $f(x)$?

Idea: Let's compare to a baseline $b \in \mathbb{R}^d$
(e.g., an average point)

Special Case: Linear Models

$$f(x) = \sum_{i=1}^d w_i x_i \quad \text{vs.} \quad f(b) = \sum_{i=1}^d w_i b_i$$

$$f(x) - f(b) = \sum_{i=1}^d w_i (x_i - b_i)$$

A: Feature i changed output by

$$\phi_i = w_i (x_i - b_i)$$

Motivation: Use ϕ_i to...

- Check for concerning behavior
- Understand model
- Select important features

General Case: Non-linear models

Challenge: Complicated interactions b/wn features

E.g., Suppose f predicts plant growth.

If temp is high, precip helps.

If temp is low, precip hurts.

Solution: Evaluate effect of feature i
in context of other features

Value Function

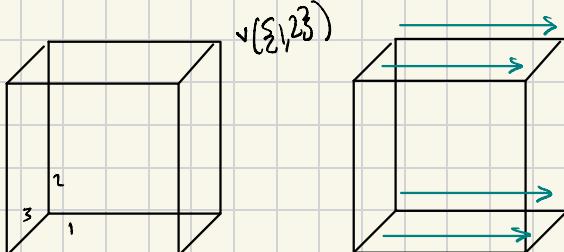
$$S \subseteq \{1, \dots, d\}$$

$x^S \in \mathbb{R}^d$ combines x, b

$$x_i^S = \begin{cases} x_i & \text{if } i \in S \\ b_i & \text{if } i \notin S \end{cases}$$

$$v(S) = f(x^S)$$

Goal: Analyze $\sum_S [v(S \cup i) - v(S)]$



Shapley Values

Axioms:

↳ Null

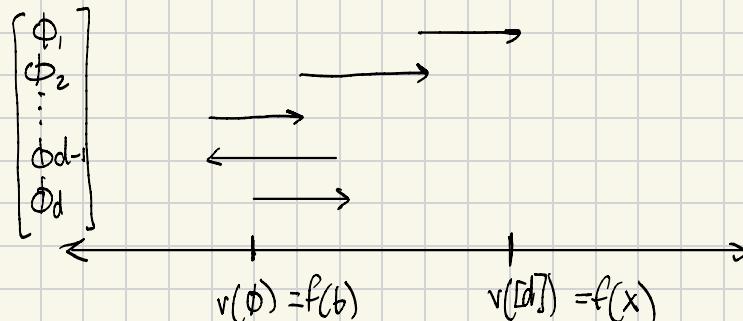
↳ Symmetry

↳ Linearity

↳ Efficiency

$$\sum_{i=1}^d \phi_i = v(\{d\}) - v(\emptyset)$$

$$\phi_i = \sum_{S \subseteq \{1, \dots, d\} \setminus \{i\}} [v(S \cup i) - v(S)] \frac{1}{d(d-1)}$$

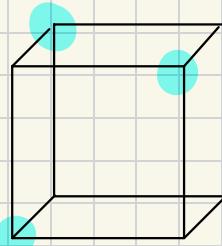


Decompose change in prediction in terms of features

Estimating Shapley Values

Challenge: 2^d subsets to evaluate

Monte Carlo Estimator approximates sum by computing several terms



$$\text{Choose } P_S = \frac{1}{d \binom{d-1}{|S|-1}} \\ \tilde{\Phi}_i^{MC} = \frac{1}{m} \sum_{j=1}^m [v(S_j \cup i) - v(S_j)] \frac{1}{d \binom{d-1}{|S_j|-1} P_{S_j}} = \frac{1}{m} \sum_{j=1}^m v(S_j \cup i) - v(S_j)$$

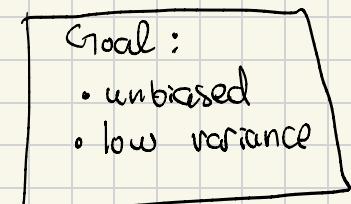
$$\mathbb{E}[\tilde{\Phi}_i^{MC}] = \mathbb{E} \left[\frac{1}{m} \sum_{j=1}^m \sum_{S \subseteq [d] \setminus i} [v(S \cup i) - v(S)] \mathbb{I}[S_j = S] \right] \\ = \dots$$

$$= \Phi_i$$

Sample m subsets instead of all 2^d

$$S_1, \dots, S_m \subseteq [d] \setminus \{i\}$$

S is sampled w/o replacement



Variance Bound

- Scalar $\text{Var}(cX) = c^2 \text{Var}(X)$

Fact: Variance without replacement \leq variance with replacement

So, $\text{Var}(\phi_i^{MC})$ at most variance when sampled with replacement

$$\begin{aligned}\text{Var}(\phi_i^{MC}) &= \text{Var} \left(\frac{1}{m} \sum_{j=1}^m \sum_{S \subseteq [d] \setminus i} [\nu(s_{v(i)}) - \nu(S)] \mathbb{I}[S_j = S] \right) \\ &\stackrel{\uparrow \text{ w/ replacement}}{=} \frac{1}{m^2} \sum_{j=1}^m \sum_{S \subseteq [d] \setminus i} [\nu(s_{v(i)}) - \nu(S)]^2 \text{Var}(\mathbb{I}[S_j = S]) \quad \text{by indep} \\ &\leq \frac{1}{m} \sum_{S \subseteq [d] \setminus i} [\nu(s_{v(i)}) - \nu(S)]^2 p_{ISI} \quad \text{by below inequality}\end{aligned}$$

$$\begin{aligned}\text{Var}(\mathbb{I}[S_j = S]) &= \mathbb{E}[\mathbb{I}[S_j = S]^2] - \mathbb{E}[\mathbb{I}[S_j = S]]^2 \\ &\leq \mathbb{E}[\mathbb{I}[S_j = S]^2] \\ &= 1 \cdot \Pr(S_j = S) + O[1 - \Pr(S_j = S)] \\ &\approx \Pr(S_j = S) = p_{ISI}\end{aligned}$$

Chebyshev's

$$\Pr(|\phi_i^{MC} - \phi_i| \geq \epsilon) \leq \frac{\text{Var}(\phi_i^{MC})}{\epsilon^2}$$

Maximum Sample Reuse

All this for only one $\hat{\phi}_i^{MC}$? Let's reuse samples!

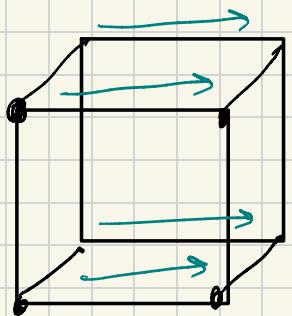
$$\begin{aligned}\hat{\phi}_i &= \sum_{S \subseteq [d] \setminus \{i\}} [v(S \cup i) - v(S)] \frac{1}{d \binom{d-1}{|S|}} \\ &= \sum_{S \subseteq [d]: i \in S} v(S) \frac{1}{d \binom{d-1}{|S|-1}} - \sum_{S \subseteq [d]: i \notin S} v(S) \frac{1}{d \binom{d-1}{|S|}} \\ &= \sum_{S \subseteq [d]} v(S) \left[\frac{\mathbb{I}[i \in S]}{d \binom{d-1}{|S|-1}} - \frac{\mathbb{I}[i \notin S]}{d \binom{d-1}{|S|}} \right]\end{aligned}$$

Estimate this! But variance will depend on $[v(S)]^2$ rather than $[v(S \cup i) - v(S)]^2$..

Last Lecture

- ↳ Midterm 11/25
- ↳ Important class on 12/2
 - proposal due 10am
 - discuss proposal in class
 - class eval

Shapley Values



• prediction given
features in S

$$v(S)$$

$$\phi_i = \sum_{S \subseteq [d] \setminus \{i\}} \frac{v(S \cup i) - v(S)}{d \binom{d-1}{|S|}}$$

= average effect of feature i
on prediction

Challenge: 2^d subsets!!

Monte Carlo

Sample $s_1, \dots, s_m \subseteq [d] \setminus \{i\}$ w.p $P_{is_i} = \frac{1}{d \binom{d-1}{|S|}}$

$$\begin{aligned}\hat{\phi}_i^{MC} &= \frac{1}{m} \sum_{S \subseteq [d] \setminus \{i\}} \frac{v(s_{ui}) - v(s)}{d \binom{d-1}{|S|} P_{is_i}} \mathbb{I}[S \text{ sampled}] \\ &= \frac{1}{m} \sum_{S \subseteq [d] \setminus \{i\}} (v(s_{ui}) - v(s)) \mathbb{I}[S \text{ sampled}]\end{aligned}$$

Goal: Let's compute variance!

Variance

$$\begin{aligned}\text{Var}(X+Y) &= \mathbb{E}[(X+Y - \mathbb{E}[X] - \mathbb{E}[Y])^2] \\ &= \mathbb{E}\left[\left(X - \mathbb{E}[X]\right)^2\right] + \mathbb{E}\left[\left(Y - \mathbb{E}[Y]\right)^2\right] + 2\mathbb{E}\left[\left(X - \mathbb{E}[X]\right)\left(Y - \mathbb{E}[Y]\right)\right] \\ &\quad \text{Var}(X) \qquad \qquad \text{Var}(Y) \qquad \qquad \text{Cov}(X, Y)\end{aligned}$$

If X, Y independent...

Variance of Monte Carlo

$$\hat{\phi}_i^{MC} = \frac{1}{m} \sum_{S \subseteq [d] \setminus \{i\}} (\nu(S \cup i) - \nu(S)) \quad \mathbb{I}[S \text{ sampled}]$$

$$= \frac{1}{m} \sum_S x_S \mathbb{I}_{S \setminus i}$$

$$\begin{aligned} \text{Var}(\hat{\phi}_i^{MC}) &= \frac{1}{m^2} \sum_S \text{Var}(x_S \mathbb{I}_{S \setminus i}) + \frac{1}{m^2} \sum_{S, S'} \text{cov}(x_S \mathbb{I}_{S \setminus i}, x_{S'} \mathbb{I}_{S' \setminus i}) \\ &= \frac{1}{m^2} \sum_S \mathbb{E}[(x_S \mathbb{I}_{S \setminus i} - \mathbb{E}[x_S \mathbb{I}_{S \setminus i}])^2] + \frac{2}{m^2} \sum_{S, S'} \mathbb{E}[(x_S \mathbb{I}_{S \setminus i} - \mathbb{E}[x_S \mathbb{I}_{S \setminus i}])(x_{S'} \mathbb{I}_{S' \setminus i} - \mathbb{E}[x_{S'} \mathbb{I}_{S' \setminus i}])] \\ &= \frac{1}{m^2} \sum_S x_S^2 (\mathbb{E}[\mathbb{I}_{S \setminus i}^2] - \mathbb{E}[\mathbb{I}_{S \setminus i}]^2) + \frac{2}{m^2} \sum_{S, S'} x_S x_{S'} \mathbb{E}[(\mathbb{I}_{S \setminus i} - \mathbb{E}[\mathbb{I}_{S \setminus i}]) (\mathbb{I}_{S' \setminus i} - \mathbb{E}[\mathbb{I}_{S' \setminus i}])] \end{aligned}$$

$$\begin{aligned} \mathbb{E}[(\mathbb{I}_{S \setminus i} - \mathbb{E}[\mathbb{I}_{S \setminus i}]) (\mathbb{I}_{S' \setminus i} - \mathbb{E}[\mathbb{I}_{S' \setminus i}])] &= \mathbb{E}[\mathbb{I}_{S \setminus i} \mathbb{I}_{S' \setminus i}] - 2 \mathbb{E}[\mathbb{I}_{S \setminus i}] \mathbb{E}[\mathbb{I}_{S' \setminus i}] + \mathbb{E}[\mathbb{I}_{S \setminus i}] \mathbb{E}[\mathbb{I}_{S' \setminus i}] \\ &= \Pr(\text{both } S, S' \text{ selected}) - \Pr(S \text{ selected}) \Pr(S' \text{ selected}) \leq 0 \end{aligned}$$

$$\text{Var}(\hat{\phi}_i^{MC}) \leq \frac{1}{m^2} \sum_{S \subseteq [d] \setminus \{i\}} [\nu(S \cup i) - \nu(S)]^2 P_{i,S}$$

$$\mathbb{E}[\mathbb{I}_{S \setminus i}^2] - \mathbb{E}[\mathbb{I}_{S \setminus i}]^2 \leq \mathbb{E}[\mathbb{I}_{S \setminus i}]$$

$$= \mathbb{E}[\mathbb{I}_{S \setminus i}]$$

$$= 1 \cdot \Pr(S) + 0 \cdot (1 - \Pr(S))$$

$$= \Pr(S) = P_{i,S}$$

Maximum Sample Reuse

Motivation: We compute $v(S \cup i) - v(S)$,
but only use it to estimate ϕ_i

Unbiased? ✓

∴

Variance? depends on $[v(S)]^2$ ∵

$$\begin{aligned}\phi_i &= \sum_{S \subseteq [d] \setminus \{i\}} \frac{v(S \cup i) - v(S)}{d \binom{d-1}{|S|-1}} \\ &= \sum_{S: i \in S} \frac{v(S)}{d \binom{d-1}{|S|-1}} - \sum_{S: i \notin S} \frac{v(S)}{d \binom{d-1}{|S|-1}} \\ &= \sum_S v(S) \left[\frac{\mathbb{I}[i \in S]}{d \binom{d-1}{|S|-1}} - \frac{\mathbb{I}[i \notin S]}{d \binom{d-1}{|S|-1}} \right]\end{aligned}$$

Sample $S_1, \dots, S_m \subseteq [d]$, use for all estimates

$$\hat{\phi}_i^{MSR} = \frac{1}{m} \sum_S v(S) \left[\frac{\mathbb{I}[i \in S]}{d \binom{d-1}{|S|-1}} - \frac{\mathbb{I}[i \notin S]}{d \binom{d-1}{|S|-1}} \right] \frac{\mathbb{I}[S \text{ sampled}]}{P_{iS}}$$

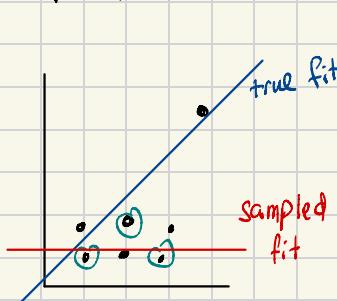
Active Learning

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

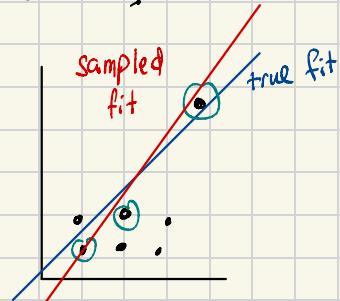
$$\begin{aligned}\Phi = \begin{bmatrix} \Phi_1 \\ \Phi_2 \\ \vdots \\ \Phi_d \end{bmatrix} &= \operatorname{argmin}_{\omega: \langle \omega, 1 \rangle = v(d) - v(\emptyset)} \|X\omega - y\|_2^2 \\ &= \operatorname{argmin}_{\omega: \langle \omega, 1 \rangle = v(d) - v(\emptyset)} \sum_{S} (v(S) - \sum_{i:i \in S} w_i)^2\end{aligned}$$

$$\begin{bmatrix} X \\ \vdots \\ X \end{bmatrix} \begin{bmatrix} w \\ \vdots \\ w^* \end{bmatrix} \approx \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} \Rightarrow \text{sample rows} \quad \begin{bmatrix} \Pi X \\ \vdots \\ \Pi X \end{bmatrix} \begin{bmatrix} w \\ \vdots \\ \hat{w} \end{bmatrix} \approx \begin{bmatrix} \Pi y \\ \vdots \\ \Pi y \end{bmatrix}$$

Q: Which rows to choose?



sampled



sampled

Idea: Choose row by its leverage

If we sample $m = \tilde{\Theta}(\frac{d}{\epsilon^2})$ rows,

$$\|X\hat{w} - y\|_2^2 \leq (1+\epsilon) \|Xw^* - y\|_2^2$$

Regression Adjustment

Sample s_1, \dots, s_m

Fit \hat{v} to v on samples

Return $\phi_i(\hat{v}) + \hat{\phi}_i^{\text{MSR}}(v - \hat{v})$

Consistent: What is $\phi_i(\hat{v}) + \phi_i(v - \hat{v})$?