

Modern approaches to the assessment of demand sensitivity based on uplift models and neural networks

D. Ustinova, A. Kurmukova, E. Avdotin, K. Onlabek

Skoltech

Skolkovo Institute of Science and Technology

ML 2021 final project defence

March 21, 2021

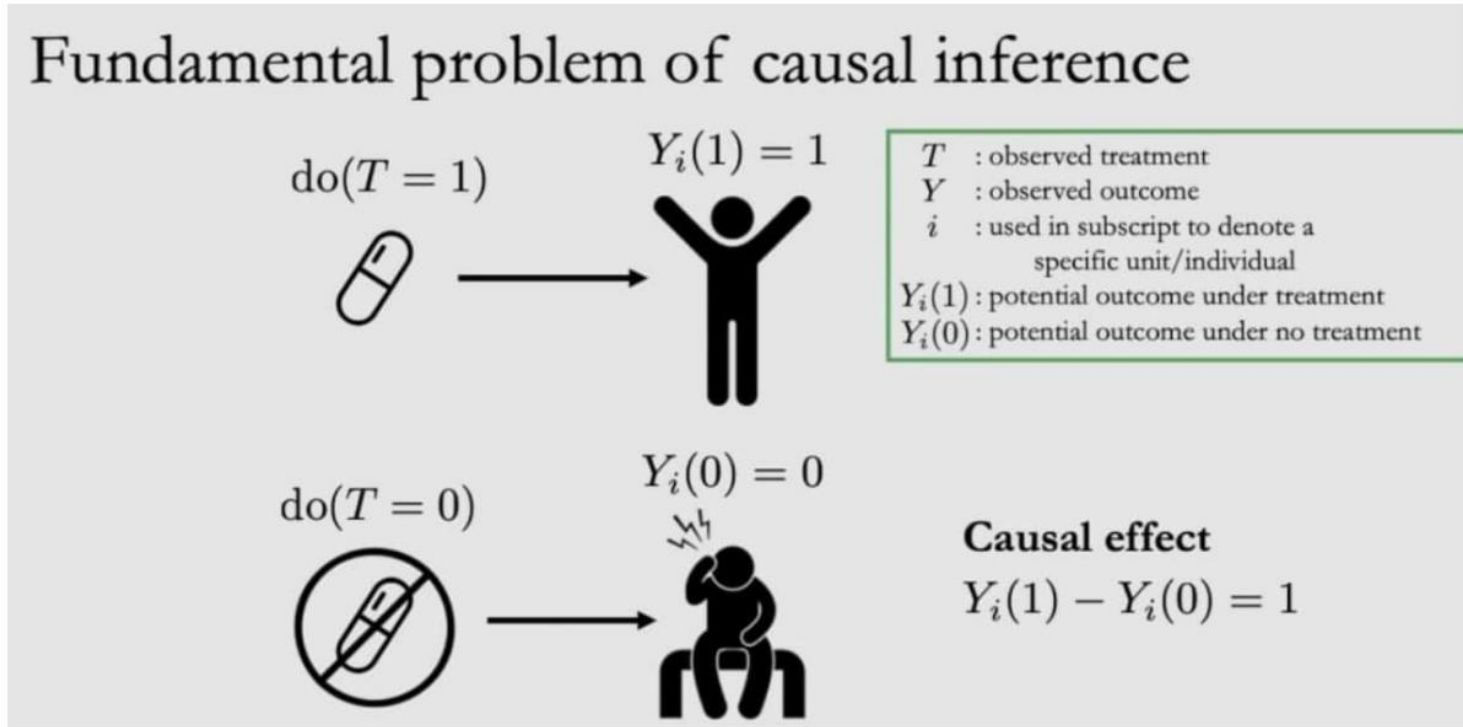
Moscow, Russia

Problem statement:

- **Uplift modeling:** predicts the incremental effect of a treatment on an individual's behaviour.
- **Example:** a company wants to send a customer an SMS with a discount offer.
- **Goal:** to detect the persuadable group of people and offer them the discount.

Response if Treated	N	Do-Not-Disturb <i>c</i>	Lost Cause <i>d</i>
	Y	Sure Thing <i>b</i>	Persuadable <i>a</i>
		Y	N
		Response if <u>not</u> treated	

Problem statement:



- Can't take and don't take the pill simultaneously
 \implies it is not possible to observe and count the uplift (or casual effect) directly.

Solution:

In classical approach:

- Casual effect:

$$\tau_i = Y_i^1 - Y_i^0.$$

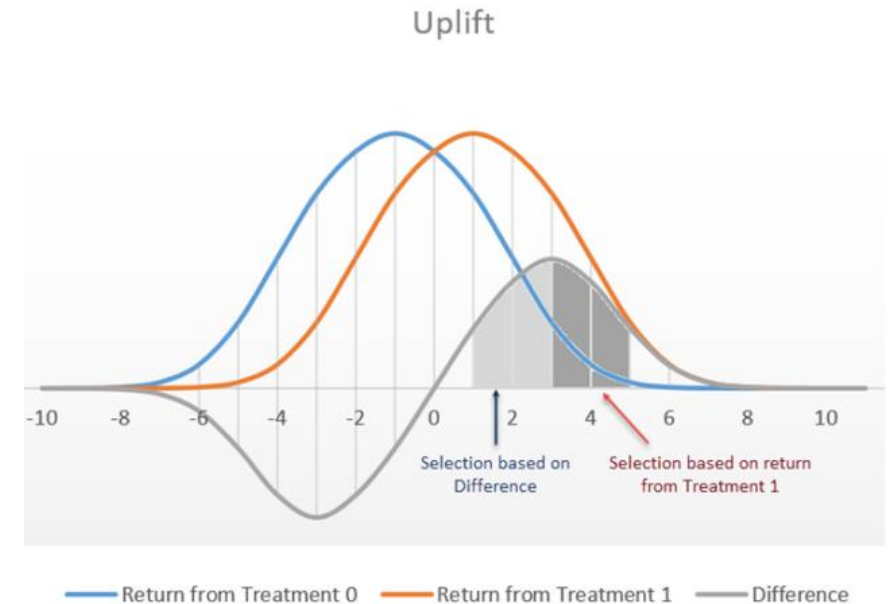
- CATE (Conditional Average Treatment Effect):

$$CATE = \mathbb{E} [Y_i^1 | X_i] - \mathbb{E} [Y_i^0 | X_i].$$

In uplift modelling:

$$\begin{aligned} uplift = \widehat{CATE} &= \mathbb{E} [Y_i | X_i = x, W_i = 1] - \\ &\quad - \mathbb{E} [Y_i | X_i = x, W_i = 0], \end{aligned}$$

where $Y_i = W_i Y_i^1 + (1 - W_i) Y_i^0$ is customer's observed response.



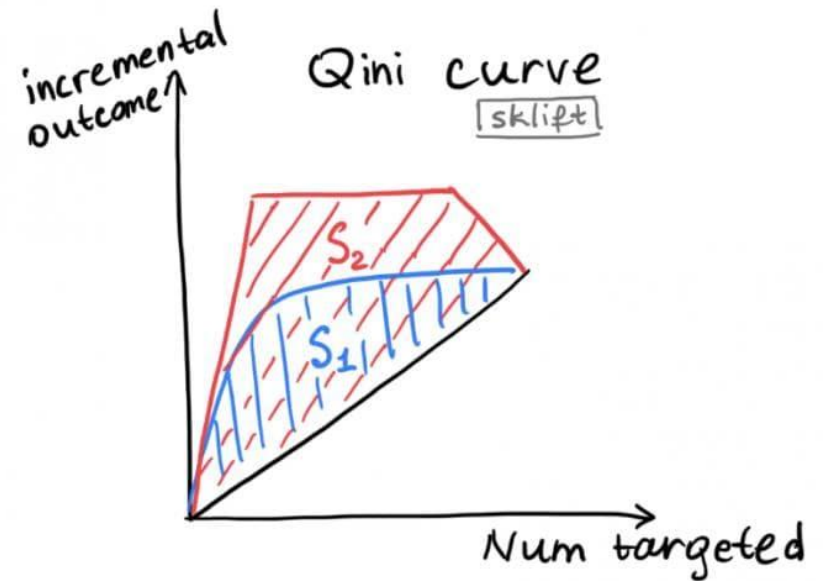
Uplift metrics

- uplift@k – the value of uplift at top k% of data
- area under uplift curve:
- area under qini curve

$$uplift_curve(t) = \left(\frac{Y_t^T}{N_t^T} - \frac{Y_t^C}{N_t^C} \right) (N_t^T + N_t^C)$$

$$qini_curve(t) = Y_t^T - \frac{Y_t^C N_t^T}{N_t^C},$$

where Y_t^T and Y_t^C are targets in treatment and control groups, respectively. N_t^T and N_t^C are sizes of groups and t is the number of targeted objects.



$$\text{Qini coefficient} = \frac{S_1}{S_2}$$

Meta-learners: S-learner

S-learner calculates the treatment effect using a single machine learning model

- fitting

$$\text{fit} \left(\begin{array}{ccc|c|c} x_{11} & \cdots & x_{1k} & w_1 & y_1 \\ \vdots & \ddots & \vdots & \cdots & \cdots \\ x_{n1} & \cdots & x_{nk} & w_n & y_n \end{array} \right)$$

$X_{train} \quad W_{train} \quad Y_{train}$

- inference

$$\text{predict}_{proba} \left(\begin{array}{ccc|c} x_{11} & \cdots & x_{1k} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_{m1} & \cdots & x_{mk} & 1 \end{array} \right) - \text{predict}_{proba} \left(\begin{array}{ccc|c} x_{11} & \cdots & x_{1k} & 0 \\ \vdots & \ddots & \vdots & \vdots \\ x_{m1} & \cdots & x_{mk} & 0 \end{array} \right) = \begin{pmatrix} u_1 \\ \vdots \\ u_m \end{pmatrix}$$

$X_{test} \quad W_1 \quad X_{test} \quad W_0 \quad \text{uplift}$

Meta-learners: T-learner

T-learner is a Two-Model approach

- fitting

$$model^T = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pk} \end{matrix}, \begin{matrix} y_1 \\ \cdots \\ y_p \end{matrix} \right), \quad model^C = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{q1} & \cdots & x_{qk} \end{matrix}, \begin{matrix} y_1 \\ \cdots \\ y_q \end{matrix} \right)$$

X_{train_treat} Y_{train_treat} $X_{train_control}$ $Y_{train_control}$

- inference

$$\begin{matrix} model^T \\ predict \\ proba \end{matrix} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mk} \end{matrix} \right) - \begin{matrix} model^C \\ predict \\ proba \end{matrix} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mk} \end{matrix} \right) = \begin{pmatrix} u_1 \\ \vdots \\ u_m \end{pmatrix}$$

X_{test} X_{test} $uplift$

Meta-learners: X-learner

$$model^T = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pk} \end{matrix}, \begin{matrix} y_1 \\ \vdots \\ y_p \end{matrix} \right), \quad model^C = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{q1} & \cdots & x_{qk} \end{matrix}, \begin{matrix} y_1 \\ \vdots \\ y_q \end{matrix} \right)$$

X_{train_treat} Y_{train_treat} $X_{train_control}$ $Y_{train_control}$

$$\hat{Y}^C = \underset{X_{train_treat}}{model^C \underset{predict}} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pk} \end{matrix} \right); \quad \hat{Y}^T = \underset{X_{train_control}}{model^T \underset{proba}} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{q1} & \cdots & x_{qk} \end{matrix} \right)$$

$$\tilde{D}^T = Y_{train_treat} - \hat{Y}^C; \quad \tilde{D}^C = \hat{Y}^T - Y_{train_control}$$

$$model_{new}^T = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pk} \end{matrix}, \begin{matrix} \tilde{d}_1^T \\ \vdots \\ \tilde{d}_p^T \end{matrix} \right), \quad model_{new}^C = fit \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{q1} & \cdots & x_{qk} \end{matrix}, \begin{matrix} \tilde{d}_1^C \\ \vdots \\ \tilde{d}_q^C \end{matrix} \right)$$

X_{train_treat} \tilde{D}^T $X_{train_control}$ \tilde{D}^C

$$g \cdot \underset{X_{test}}{model_{new}^C \underset{predict}} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mk} \end{matrix} \right) + (1 - g) \cdot \underset{X_{test}}{model_{new}^T \underset{predict}} \left(\begin{matrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mk} \end{matrix} \right) = \begin{pmatrix} u_1 \\ \vdots \\ u_m \end{pmatrix}$$

$uplift$

Meta-learners: R-learner

R-learner is a two-step algorithm

- Fitting two base models using cross-validation and predict on holdout fold i:
 - the outcomes $m^{(-i)}(X_i)$
 - the propensity scores $e^{(-i)}(X_i)$
- Minimising of R-loss:

$$\tau^*() = \operatorname{argmin}_{\tau} \left\{ \frac{1}{n} \sum_{i=1}^n [(Y_i - \hat{m}^{(-i)}(X_i)) - (W_i - \hat{e}^{(-i)}(X_i))\tau(X_i)]^2 + \Lambda(\tau()) \right\},$$

Tree-based approaches: Uplift Tree

The tree is constructed to maximize the difference between distribution divergence for the treatment and control groups before and after the split. The best split in each node is chosen based on the gain:

$$D_{gain} = D_{aftersplit}(P^T, P^C) - D_{beforesplit}(P^T, P^C).$$

Three approaches to calculate the divergence:

- the Kullback-Leibler (KL) divergence:

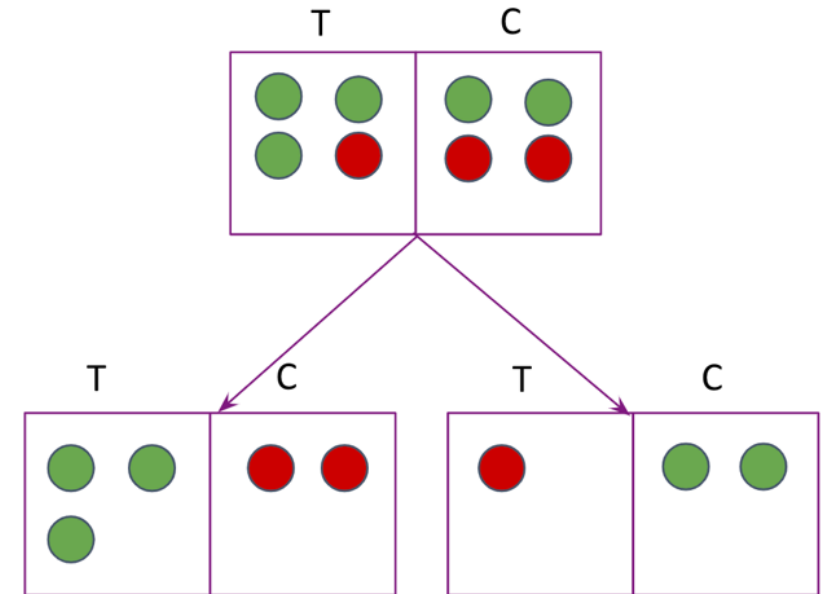
$$KL(P : Q) = \sum_i p_i \log \frac{p_i}{q_i};$$

- the Euclidean distance:

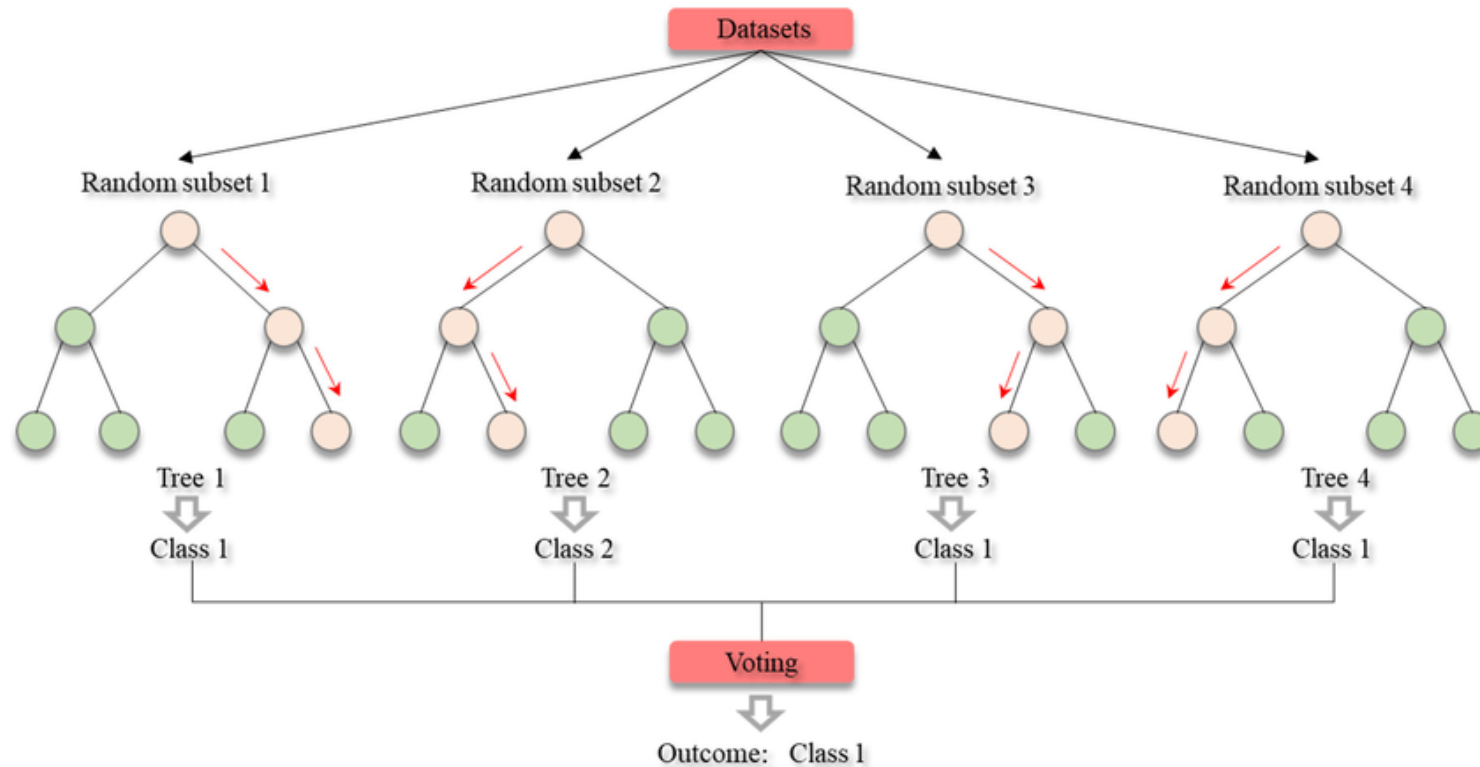
$$ED(P : Q) = \sum_i (p_i - q_i)^2;$$

- the Chi-squared divergence:

$$\xi^2(P : Q) = \sum_i \frac{(p_i - q_i)^2}{q_i}.$$



Tree-based approaches: Random Forest



Boosting

In the project we considered Uplift AdaBoost algorithm:

1. Initialize weights $w_{1,i}^T, w_{1,i}^C$
2. For $m \leftarrow 1, \dots, M$
 - (a) $w_{m,i}^T \leftarrow \frac{w_{m,i}^T}{\sum_j w_{m,j}^T + \sum_j w_{m,j}^C}; w_{m,i}^C \leftarrow \frac{w_{m,i}^C}{\sum_j w_{m,j}^T + \sum_j w_{m,j}^C}$
 - (b) Build a base model h_m on \mathcal{D} with $w_{m,i}^T, w_{m,i}^C$
 - (c) Compute the treatment and control errors $\epsilon_m^T, \epsilon_m^C$
 - (d) Compute $\beta_m^T(\epsilon_m^T, \epsilon_m^C), \beta_m^C(\epsilon_m^T, \epsilon_m^C)$
 - (e) If $\beta_m^T = \beta_m^C = 1$ or $\epsilon_m^T \notin (0, \frac{1}{2})$ or $\epsilon_m^C \notin (0, \frac{1}{2})$:
 - i. choose random weights $w_{m,i}^T, w_{m,i}^C$
 - ii. continue with next boosting iteration
 - (f) $w_{m+1,i}^T \leftarrow w_{m,i}^T \cdot (\beta_m^T)^{1[h_m(x_i^T)=y_i^T]}$
 - (g) $w_{m+1,i}^C \leftarrow w_{m,i}^C \cdot (\beta_m^C)^{1[h_m(x_i^C)=1-y_i^C]}$
 - (h) $\beta_m \leftarrow \min\{\beta_m^T, \beta_m^C\}$
 - (i) Add h_m with coefficient β_m to the ensemble

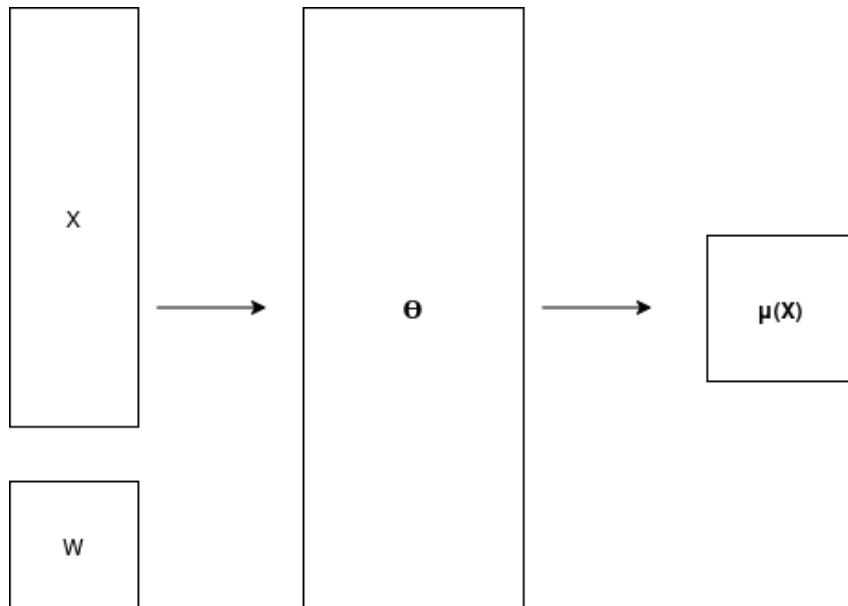
Output: The final hypothesis

$$h_f(x) = \begin{cases} 1 & \text{if } \sum_{m=1}^M \left(\log \frac{1}{\beta_m} \right) h_m(x) \geq \frac{1}{2} \sum_{m=1}^M \log \frac{1}{\beta_m}, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Neural Network for uplift modeling

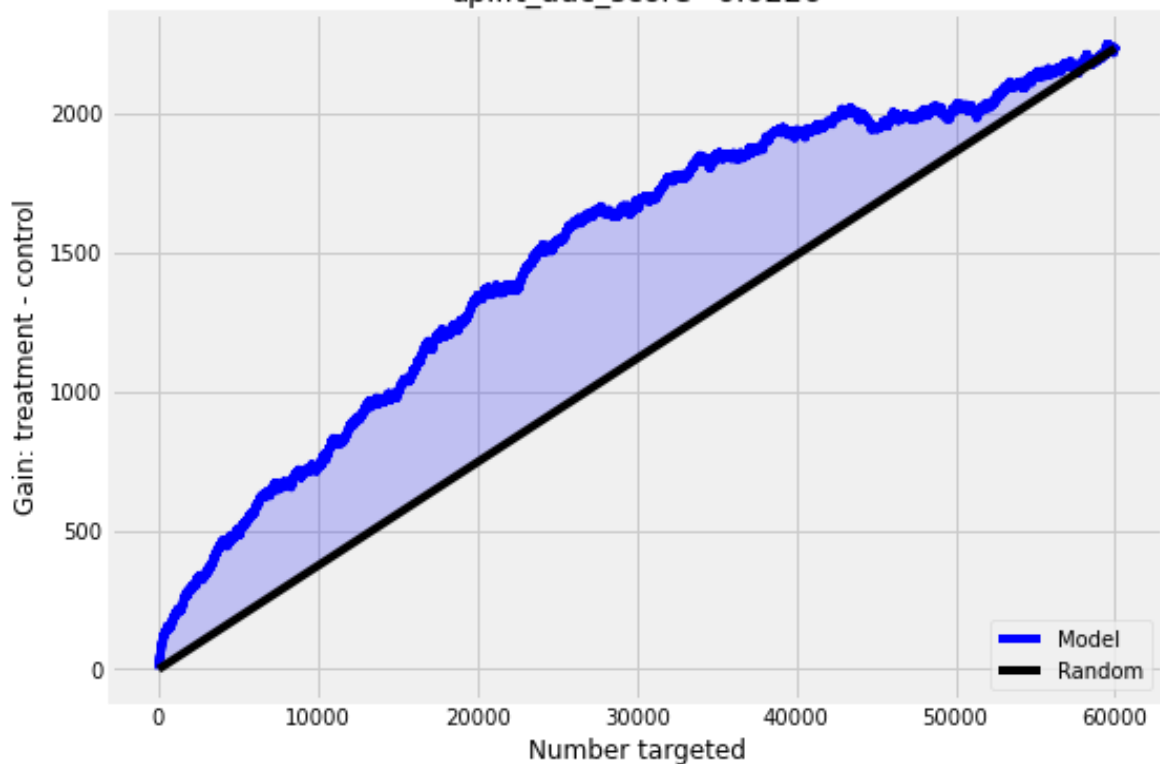
- Transformed outcome is used as target
- MLP architecture
- Transformed output depending on treatment
- Training and validation on X5 dataset

$$\hat{y}_i = \begin{cases} \mu(x_i), & \text{if } w_i = 1 \\ -\mu(x_i), & \text{if } w_i = 0 \end{cases}$$



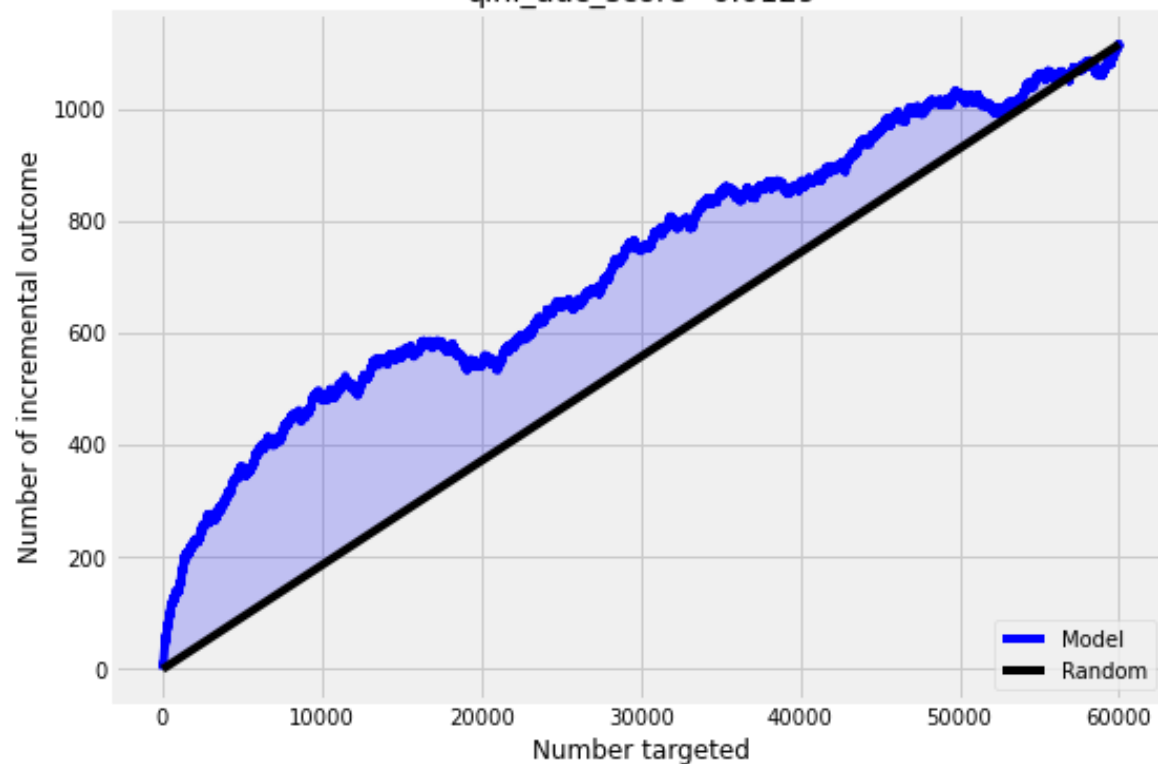
Results: meta-learners

Uplift curve
uplift_auc_score=0.0226



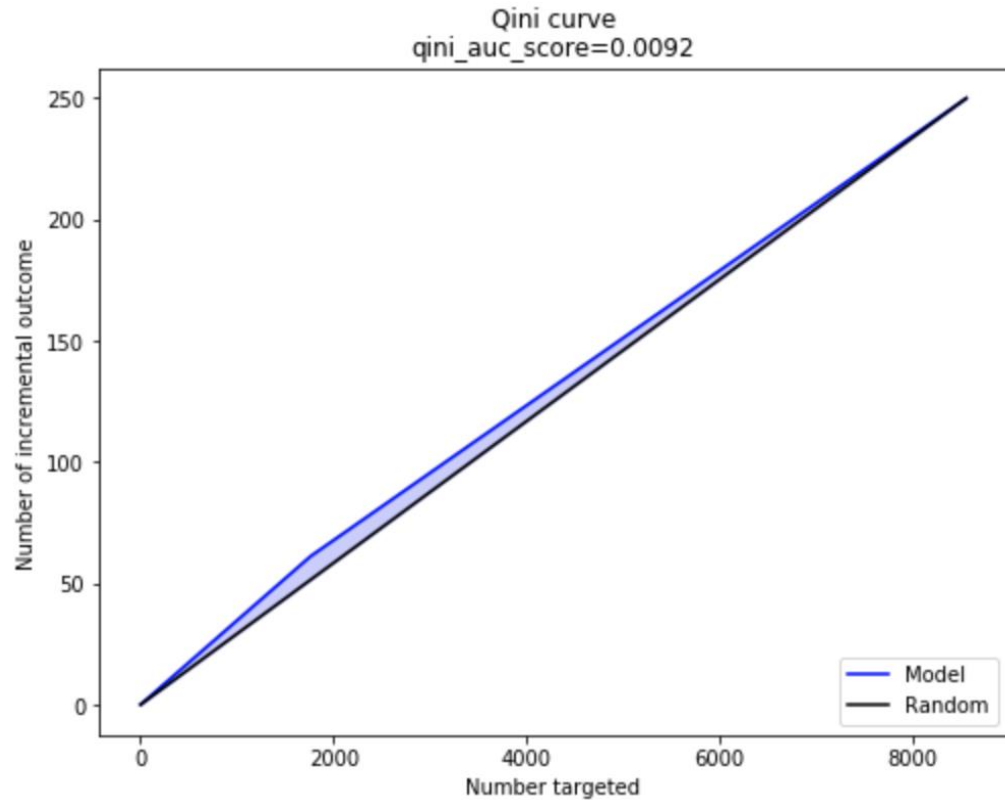
Uplift curve for S learner
on X5 dataset

Qini curve
qini_auc_score=0.0129

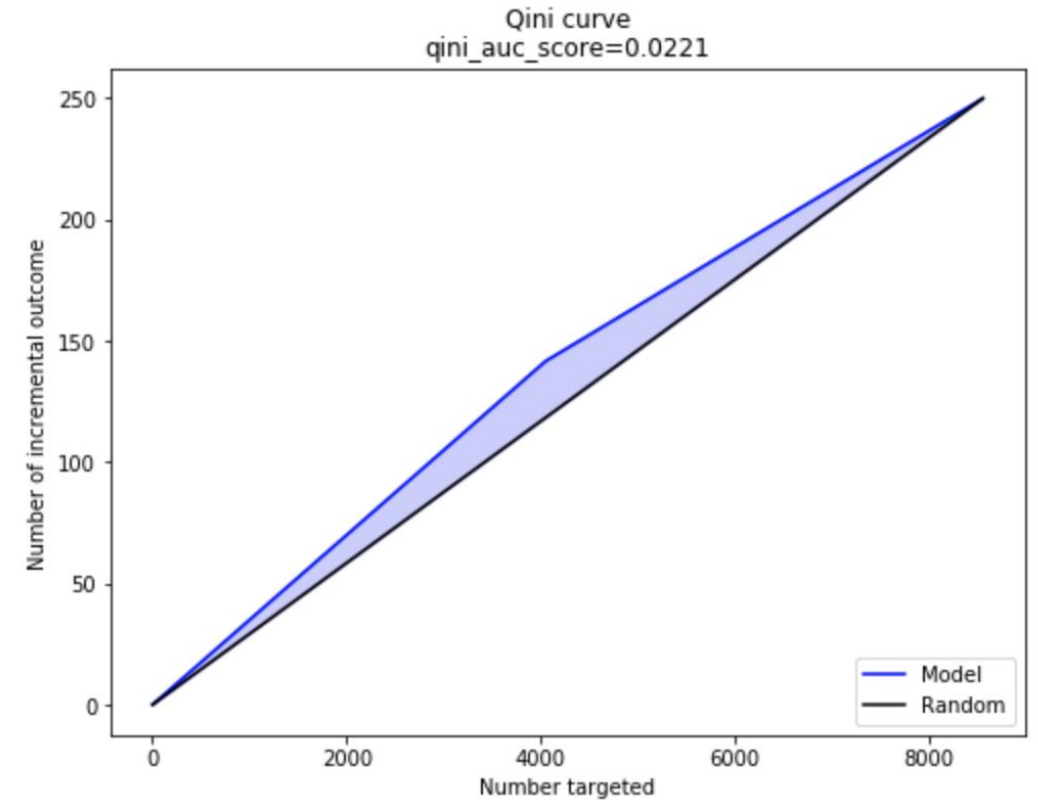


Qini curve for T learner on
X5 dataset

Results: random forest and boosting

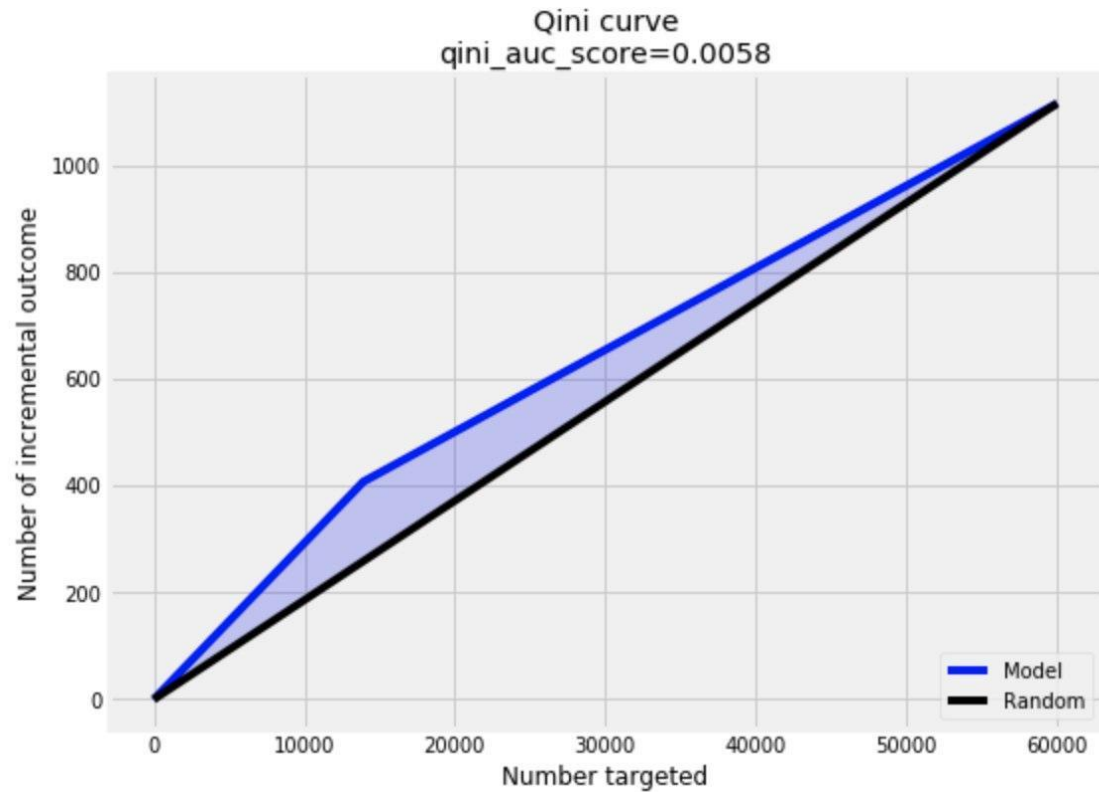


Qini curve for random forest
On MineThatData

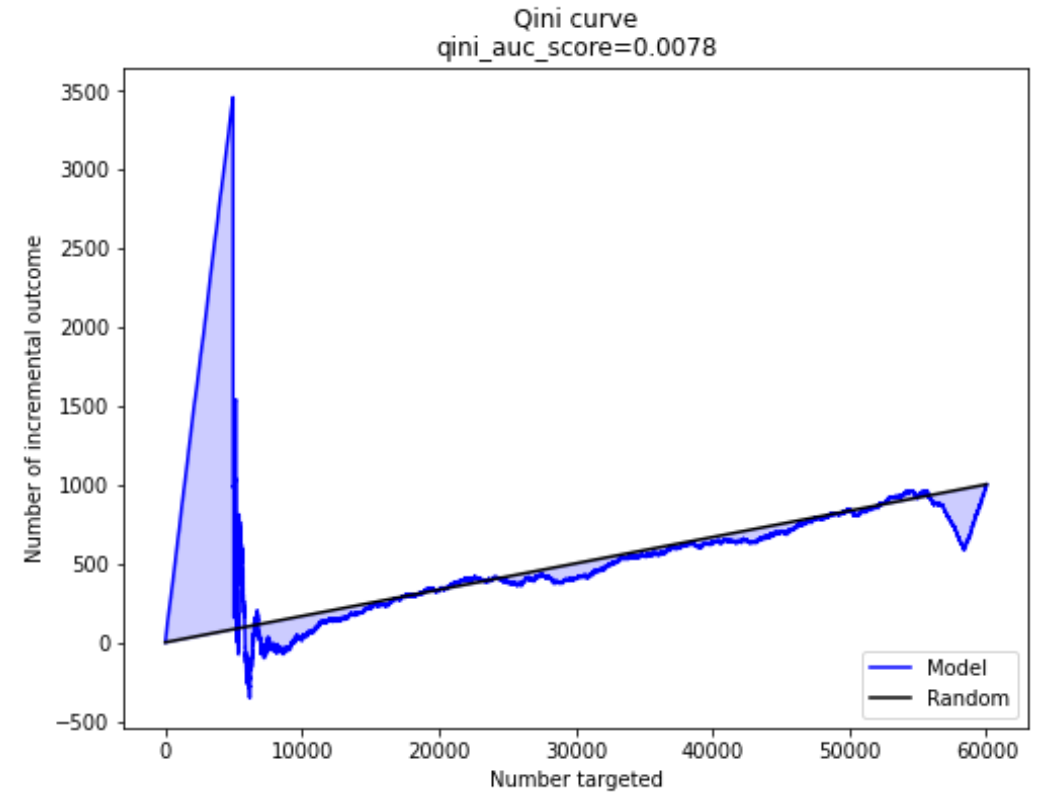


Qini curve for boosting
On MineThatData

Results: trees and neural network



Qini curve for tree with
depth 8, X5 dataset



Qini curve for
Neural Network, X5 dataset

Conclusion

- Meta-learners (S-, T-, X- and R-learners)
- Tree-based approaches (Uplift Tree, Random Forest)
- Boosting (AdaBoost)
- Neural networks
- Comparison on 3 datasets

Thank you for your
attention!