# 1   Variable Selection

Let $y$ be response variable and $x$ be explanatory variables or covariates. Given i.i.d. samples $(x, y) \in \mathbb{R}^p \times \mathbb{R}$ from the joint distribution $p_{x,y}$, we are interested in asking the question

*which of the many covariates $x_1, \cdots, x_p$ does the response $y$ depend on?*

assuming that the response does depend on a sparse set of variables. In reality, we are interested in the causal relationship. However, quantifying causal effects requires interventions and not possible from purely observational data. A natural relaxation is to find covariates dependent (in a statistical sense) on the response, conditioned on all other observed features [1]. Formally, we want to find smallest $\mathcal{S} \subset [p]$ s.t.

$$y \perp\!\!\!\perp x_{\mathcal{S}} \mid x_{\backslash \mathcal{S}}$$

A natural interpretation is that the other variables $x_{\backslash \mathcal{S}}$ do not provide additional information about $y$. If we think of $\mathcal{G}$ as graph representing the joint distribution $p_{x,y}$, then $\mathcal{S}$ is the markov blanket for node $y$. We can pose the problem of finding the Markov blanket of $y$ as



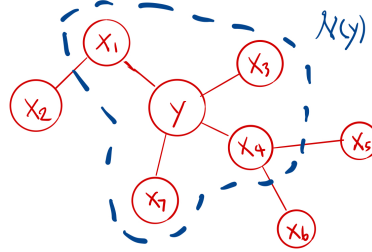Figure 1: $\mathcal{S} = \{x_1, x_3, x_4, x_7\}$

a multiple binary hypothesis test

$$H_0^{(j)} : y \perp\!\!\!\perp x_j \mid x_{\backslash \{j\}} \qquad \text{for} \quad j = 1, \cdots, p$$

Let $\mathcal{H}_0 = \left\{ x_j \mid H_0^{(j)} \text{ holds} \right\}$ be the set of truly irrelevant covariates. In essence we are interested in maximizing *power* while controlling the number of false positives. A global threshold for p-values of each tests is overly conservative for large $p$, an alternative approach is to control *false discovery rate* (FDR) [2].

$$\text{maximize}_{\hat{\mathcal{S}} \subset [p]} \quad \mathbb{E}\left[ \frac{\hat{\mathcal{S}} \setminus \mathcal{H}_0}{|\hat{\mathcal{S}}|} \right] \qquad \text{(maximize power)}$$

$$\text{subject to} \quad \mathbb{E}\left[ \frac{\hat{\mathcal{S}} \cap \mathcal{H}_0}{|\hat{\mathcal{S}}|} \right] \leq q \qquad \text{(control FDR)}$$

# 2   Model-X Knockoff

# References

[1] Jaime Roquero Gimenez, Amirata Ghorbani, and James Zou. "Knockoffs for the mass: new feature importance statistics with false discovery guarantees". In: *arXiv:1807.06214 [cs, stat]* (May 28, 2019). arXiv: 1807.06214. URL: http://arxiv.org/abs/1807.06214 (visited on 04/17/2020).

[2] Yoav Benjamini and Yosef Hochberg. "Controlling The False Discovery Rate - A Practical And Powerful Approach To Multiple Testing". In: *J. Royal Statist. Soc., Series B* 57 (Nov. 30, 1995), pp. 289–300. DOI: 10.2307/2346101.