

A. Theoretical Analysis of Local Feature FDR

Definition A.1. We call a mapping $\sigma : \mathbb{R}^d \rightarrow \mathcal{P}([d])$ a generic swap or a swap. In addition, we say that a swap is a local swap if for any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$,

$$\mathbf{x}_{[d] \setminus \sigma(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma(\mathbf{x})} \Rightarrow \sigma(\mathbf{x}) = \sigma(\mathbf{z})$$

Given a mapping

$$(\mathbf{F}, \tilde{\mathbf{F}}) : \begin{cases} \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d \\ (\mathbf{x}, \tilde{\mathbf{x}}) \mapsto (F_1(\mathbf{x}, \tilde{\mathbf{x}}), \dots, F_d(\mathbf{x}, \tilde{\mathbf{x}}), \\ \tilde{F}_1(\mathbf{x}, \tilde{\mathbf{x}}), \dots, \tilde{F}_d(\mathbf{x}, \tilde{\mathbf{x}})) \end{cases}$$

and swap σ , define the operation $[\mathbf{F}, \tilde{\mathbf{F}}]_{\text{swap}(\sigma)}$ as the mapping $[\mathbf{F}, \tilde{\mathbf{F}}]_{\text{swap}(\sigma)} : (\mathbf{x}, \tilde{\mathbf{x}}) \mapsto [(\mathbf{F}, \tilde{\mathbf{F}})(\mathbf{x}, \tilde{\mathbf{x}})]_{\text{swap}(\sigma(\mathbf{x}))}$.

It is important to clearly identify the input space of the swap σ . The output of the swap operation on a mapping is again a mapping with the same input space, so for example we can iterate swap operations on a given mapping $(\mathbf{F}, \tilde{\mathbf{F}})$: for σ_1, σ_2 swaps the following is well defined: $[[\mathbf{F}, \tilde{\mathbf{F}}]_{\text{swap}(\sigma_1)}]_{\text{swap}(\sigma_2)}$. Both σ_1, σ_2 are evaluated on a same given point \mathbf{x} of the space, and the final mapping corresponds to a point-wise concatenation of the swaps. As a consequence, the order of the swap operations does not matter. In particular, for $\sigma = \sigma_1 = \sigma_2$, we have

$$[[\mathbf{F}, \tilde{\mathbf{F}}]_{\text{swap}(\sigma)}]_{\text{swap}(\sigma)} = (\mathbf{F}, \tilde{\mathbf{F}})$$

Such mappings from $\mathbb{R}^d \times \mathbb{R}^d$ to $\mathbb{R}^d \times \mathbb{R}^d$ can be concatenated: for $(\mathbf{F}, \tilde{\mathbf{F}})$, and $(\mathbf{G}, \tilde{\mathbf{G}})$ we denote $(\mathbf{F}, \tilde{\mathbf{F}}) \circ (\mathbf{G}, \tilde{\mathbf{G}})$ the concatenated mapping. Note that concatenation and swap operations do not behave well. For example, in the general case,

$$\begin{aligned} [(\mathbf{F}, \tilde{\mathbf{F}}) \circ (\mathbf{G}, \tilde{\mathbf{G}})]_{\text{swap}(\sigma)} &\neq (\mathbf{F}, \tilde{\mathbf{F}}) \circ [(\mathbf{G}, \tilde{\mathbf{G}})]_{\text{swap}(\sigma)} \\ [(\mathbf{F}, \tilde{\mathbf{F}}) \circ (\mathbf{G}, \tilde{\mathbf{G}})]_{\text{swap}(\sigma)} &\neq [(\mathbf{F}, \tilde{\mathbf{F}})]_{\text{swap}(\sigma)} \circ (\mathbf{G}, \tilde{\mathbf{G}}) \end{aligned}$$

One particular case is the identity mapping that we denote $(\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id})$, i.e. $\mathbf{F}_j^{Id}(\mathbf{x}, \tilde{\mathbf{x}}) = x_j$ and $\tilde{\mathbf{F}}_j^{Id}(\mathbf{x}, \tilde{\mathbf{x}}) = \tilde{x}_j$. Whenever we consider random variables $\mathbf{X}, \tilde{\mathbf{X}}$ we denote $[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)} := [(\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id})]_{\text{swap}(\sigma)}(\mathbf{X}, \tilde{\mathbf{X}})$ the random variable resulting from applying the swapped identity map. In addition, if σ is constant equal to $S \subset [d]$, then we go back to the previous definition of swap $[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}$ in (Candès et al., 2018). Whenever we consider an identity mapping and a local swap σ , we have the immediate following result:

$$(\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}) = [(\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id})]_{\text{swap}(\sigma)} \circ [(\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id})]_{\text{swap}(\sigma)}$$

Our goal now is to show that the exchangeability condition that defines a knockoff variable implies a stronger distributional result.

Proposition A.1. Let σ be a local swap. If $\tilde{\mathbf{X}}$ is a knockoff random variable for \mathbf{X} (i.e. satisfies exchangeability), then

$$[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)} \stackrel{d}{=} [\mathbf{X}, \tilde{\mathbf{X}}] \quad (8)$$

which we refer to as local exchangeability. If $\sigma \subset \mathcal{H}_0^0$, then

$$[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y \stackrel{d}{=} [\mathbf{X}, \tilde{\mathbf{X}}], Y \quad (9)$$

We prove this result in Section B.1. This result extends the exchangeability property and the Lemma 3.2 in (Candès et al., 2018). Instead of swapping a fixed set of features, we now allow the swapping indices to depend on the features. Notice that the knockoffs $\tilde{\mathbf{X}}$ are constructed as in the general case, the local exchangeability does not require a different definition for knockoff variables. We extend the definition of a swap to probability distributions: for $\mu \in Pr(\mathbb{R}^d \times \mathbb{R}^d)$, we denote $\mu_{\text{swap}(\sigma)} := \mathcal{L}([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)})$ whenever $\mu = \mathcal{L}([\mathbf{X}, \tilde{\mathbf{X}}])$. Abusing notation, whenever $\mu = \mathcal{L}([\mathbf{X}, \tilde{\mathbf{X}}], Y)$ we will still denote $\mu_{\text{swap}(\sigma)} := \mathcal{L}([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y)$.

Local Feature Statistics The next step is to extend the construction of feature statistics to the local setting.

Definition A.2. Define local importance scores as a mapping:

$$\Phi : \begin{cases} Pr(\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}) \rightarrow (\mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d) \\ \mu \mapsto (\mathbf{T}_\mu, \tilde{\mathbf{T}}_\mu) \end{cases} \quad (10)$$

where

$$(\mathbf{T}_\mu, \tilde{\mathbf{T}}_\mu) : \begin{cases} \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d \\ \mathbf{z} \mapsto (\mathbf{T}_\mu(\mathbf{z}), \tilde{\mathbf{T}}_\mu(\mathbf{z})) \end{cases} \quad (11)$$

such that, for any $S \subset [d]$, we have

$$\Phi(\mu_{\text{swap}(S)}) = [\Phi(\mu)]_{\text{swap}(S)}$$

For $r > 0$, we say that such importance scores Φ are r -local if, for any $\mu \in Pr(\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})$, we have that $\Phi(\mu)(\mathbf{z}) = (\mathbf{T}_\mu(\mathbf{z}), \tilde{\mathbf{T}}_\mu(\mathbf{z}))$ only depends on μ through the restriction of μ to $B(\mathbf{z}, r) \times B(\mathbf{z}, r) \times \mathbb{R}$. That is, if μ, μ' are two probability measures on $\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}$ such that they coincide on $B(\mathbf{z}, r) \times B(\mathbf{z}, r) \times \mathbb{R}$, then $(\mathbf{T}_\mu(\mathbf{z}), \tilde{\mathbf{T}}_\mu(\mathbf{z})) = (\mathbf{T}_{\mu'}(\mathbf{z}), \tilde{\mathbf{T}}_{\mu'}(\mathbf{z}))$.

The next goal consists in translating the swap operation in $\mu_{\text{swap}(\sigma)} = \mathcal{L}([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y)$ into a swap of $[\mathbf{T}_\mu, \tilde{\mathbf{T}}_\mu]_{\text{swap}(\sigma)}$. This step does not require $\tilde{\mathbf{X}}$ to be a knockoff of \mathbf{X} : in what follows we do not make any assumption on μ . Notice that the swap operation has been defined (Definition 5.1) as a transformation of a mapping $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d$, but it can be immediately extended to mappings $\mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d$. We are able to relate $[\mathbf{T}_{\mu_{\text{swap}(\sigma)}}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(\sigma)}}]$ and $[\mathbf{T}_\mu, \tilde{\mathbf{T}}_\mu]_{\text{swap}(\sigma)}$ if we assume locality constraints on the importance scores.

Definition A.3. For $\sigma : \mathbb{R}^d \rightarrow \mathcal{P}([d])$ local swap, and $r > 0$, define

$$A^r = \{z \in \mathbb{R}^d, \sigma(z) = \sigma(y) \forall y \in B(z, r)\}$$

$$\sigma^r : z \mapsto \begin{cases} \sigma(z) & \text{if } z \in A^r \\ \emptyset & \text{else} \end{cases}$$

Proposition A.2. Assume there exists $r > 0$ such that the importance scores are r -local. Then, assuming A^r is non-empty, for $z \in A^r$, we have:

$$[(\mathbf{T}_\mu(z), \tilde{\mathbf{T}}_\mu(z))]_{\text{swap}(\sigma(z))} = [\mathbf{T}_{\mu_{\text{swap}(\sigma)}(z)}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(\sigma)}(z)}]$$

We prove this result in Section B.2. Whenever $\mu = \mathcal{L}([\tilde{\mathbf{X}}, \tilde{\mathbf{X}}], Y)$ where $\tilde{\mathbf{X}}$ is a knockoff of \mathbf{X} , consider a local swap σ such that $\sigma \subset \mathcal{H}_0^0$. By proposition 5.1, we get that $\mu = \mu_{\text{swap}(\sigma)}$, and therefore $(\mathbf{T}_\mu, \tilde{\mathbf{T}}_\mu) = (\mathbf{T}_{\mu_{\text{swap}(\sigma)}}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(\sigma)}})$. In practice, we can imagine that instead of using μ as an input when constructing importance scores, we use $\hat{\mu}_n$, an empirical measure defined by the dataset of n i.i.d. samples from μ that we feed as input to the algorithm. Therefore a consequence of proposition 5.1 will be that, taking also into account the eventual randomness when generating importance scores, we have for any $z \in \mathbb{R}^d$,

$$(\mathbf{T}_{\hat{\mu}_n}(z), \tilde{\mathbf{T}}_{\hat{\mu}_n}(z)) \stackrel{d}{=} (\mathbf{T}_{\hat{\mu}_n, \text{swap}(\sigma)}(z), \tilde{\mathbf{T}}_{\hat{\mu}_n, \text{swap}(\sigma)}(z))$$

We now combine this equality with that of Proposition 5.2. Fix $r > 0$ such that we have r -local importance scores. Consider a set of L points $z_1, \dots, z_N \in \mathbb{R}^d$ that are pairwise $2r$ far apart, that is, for any $1 \leq l, l' \leq N$, $\|z_l - z_{l'}\|_\infty \geq 2r$.

Proposition A.3. σ^r is a local swap and if $\sigma \subset \mathcal{H}_0^0$, then $\sigma^r \subset \mathcal{H}_0^r$. Furthermore,

$$[\mathbf{T}_{\hat{\mu}_n}(z_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(z_l)]_{1 \leq l \leq L} \stackrel{d}{=} [[\mathbf{T}_{\hat{\mu}_n}(z_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(z_l)]_{\text{swap}(\sigma^r(z_l))}]_{1 \leq l \leq L}$$

We prove this proposition in Section B.3. This allows us to conclude. Fix a target FDR level $q \in (0, 1)$. Indeed, Proposition B.3 directly implies the flip-sign condition of Lemma 3.3 in (Candès et al., 2018). Independently for each z_l , consider an independent random variable $\epsilon_l = (\epsilon_{l,1}, \dots, \epsilon_{l,d})$, where for each $1 \leq l \leq L$, and $1 \leq j \leq d$ we have $\epsilon_j^l = 1$ if $j \notin \sigma^r(z_l)$, and a Rademacher random variable if $j \in \sigma^r(z_l)$. Then denoting $\sigma_\epsilon^r(z_l) = \sigma^r(z_l) \cap \{j : \epsilon_{l,j} = -1\}$, we have:

$$(\mathbf{T}_{\hat{\mu}_n}(z_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(z_l)) \stackrel{d}{=} [\mathbf{T}_{\hat{\mu}_n}(z_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(z_l)]_{\text{swap}(\sigma_\epsilon^r(z_l))}$$

As a consequence, denoting

$$\mathbf{W}_{\hat{\mu}_n}(z_l) = \mathbf{T}_{\hat{\mu}_n}(z_l) - \tilde{\mathbf{T}}_{\hat{\mu}_n}(z_l)$$

we get that

$$\epsilon_l \odot \mathbf{W}_{\hat{\mu}_n}(z_l) \stackrel{d}{=} \mathbf{W}_{\hat{\mu}_n}(z_l)$$

where the symbol \odot indicates component-wise multiplication. Now, given that $\|z_l - z_{l'}\|_\infty \geq 2r$, this equality holds uniformly for $1 \leq l \leq L$. The random choice of the swap $\sigma_\epsilon^r(z_l)$ is done independently of a random swap $\sigma_\epsilon^r(z_{l'})$ at another point $z_{l'}$. We conclude that the knockoff selection procedure now applies to each of these vectors in an independent way: that is, for each $1 \leq l \leq L$, setting

$$\hat{\tau}_l = \min \left\{ t > 0 : \frac{1 + \#\{j : [\mathbf{W}_{\hat{\mu}_n}(z_l)]_j \leq -t\}}{\#\{j : [\mathbf{W}_{\hat{\mu}_n}(z_l)]_j \geq t\}} \leq q \right\}$$

allows to construct selection sets $\hat{\mathcal{S}}_l = \{j : [\mathbf{W}_{\hat{\mu}_n}(z_l)]_j \geq \hat{\tau}_l\}$, that control FDR^r given that initially $\sigma^r(z) \subset \mathcal{H}_0^r(z)$.

That is, according to Theorem 3.4 in (Candès et al., 2018), the set $\hat{\mathcal{S}}_l$ is such that

$$\mathbb{E} \left[\frac{|\hat{\mathcal{S}}_l \cap \mathcal{H}_0^r(\mathbf{x})|}{1 \vee |\hat{\mathcal{S}}_l|} \right] \leq q$$

hence the result.

B. Proofs

B.1. Proof of Proposition 5.1

Proof. We begin the proof with two lemmas:

Lemma B.1. We can decompose a local swap $\sigma : \mathbb{R}^d \rightarrow \mathcal{P}([d])$ into $\sigma_i : \mathbb{R}^d \rightarrow \mathcal{P}([d])$ such that for every $i \in [d]$:

$$\begin{cases} \text{Im}(\sigma_i) = \{\emptyset, \{i\}\} \\ \sigma(\mathbf{x}) = \bigsqcup_{i=1}^d \sigma_i(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^d \\ \sigma_i \text{ is a local swap} \end{cases}$$

We will denote by $\sigma = \bigsqcup_{i=1}^d \sigma_i$ the property $\sigma(\mathbf{x}) = \bigsqcup_{i=1}^d \sigma_i(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^d$.

Proof. Define, for every $i \in [d]$,

$$\sigma_i(\mathbf{x}) = \begin{cases} \emptyset & \text{if } i \notin \sigma(\mathbf{x}) \\ \{i\} & \text{if } i \in \sigma(\mathbf{x}) \end{cases}$$

We need to show that σ_i is a local swap. Let $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ such that $\mathbf{x}_{[d] \setminus \sigma_i(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma_i(\mathbf{x})}$. If $\sigma_i(\mathbf{x}) = \emptyset$ then $\mathbf{x} = \mathbf{z}$ and therefore $\sigma_i(\mathbf{x}) = \sigma_i(\mathbf{z})$. If $\sigma_i(\mathbf{x}) = \{i\}$, then $i \in \sigma(\mathbf{x})$, so $\mathbf{x}_{[d] \setminus \sigma_i(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma_i(\mathbf{x})}$ implies $\mathbf{x}_{[d] \setminus \sigma(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma(\mathbf{x})}$ and therefore $\sigma(\mathbf{x}) = \sigma(\mathbf{z})$ given that σ is a local swap. But then we have $\sigma_i(\mathbf{x}) = \sigma_i(\mathbf{z})$ by definition of σ_i . \square

Lemma B.2. *Assume that we have a partition of a local swap σ into two local swaps σ^a, σ^b : $\sigma = \sigma^a \sqcup \sigma^b$. Then we have :*

$$[\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma)} = [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma^b)} \circ [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma^a)}$$

Proof. It is crucial to notice that, given our previous swap definition for an identity mapping, the $\text{swap}(\sigma^b)$ operator applies to the output of the swapped vector by σ^a . In order to prove the result we need to show that $\sigma^b(\mathbf{x}) = \sigma^b(\mathbf{z})$ where \mathbf{z} is the vector of the first d coordinates of $[\mathbf{x}, \tilde{\mathbf{x}}]_{\text{swap}(\sigma^a)}$. Given that, for any $\mathbf{x} \in \mathbb{R}^d$ we have $\mathbf{x}_{[d] \setminus \sigma^a(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma^a(\mathbf{x})}$ and $\sigma^a(\mathbf{x}) \subset \sigma(\mathbf{x})$, we get $\sigma(\mathbf{x}) = \sigma(\mathbf{z})$ and same equality with σ^a , as both σ, σ^a are local swaps. That implies $\sigma^b(\mathbf{x}) = \sigma^b(\mathbf{z})$ as $\sigma = \sigma^a \sqcup \sigma^b$. Notice that the order does not matter when composing the two swapped identity mappings. \square

In order to prove equation 4, we write it in terms of mappings: we want to show that if $(\mathbf{X}, \tilde{\mathbf{X}})$ satisfy exchangeability, then

$$\begin{aligned} [\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)} &= [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma)}(\mathbf{X}, \tilde{\mathbf{X}}) \\ &\stackrel{d}{=} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}](\mathbf{X}, \tilde{\mathbf{X}}) = [\mathbf{X}, \tilde{\mathbf{X}}] \end{aligned}$$

With Lemma B.1 we decompose $\sigma = \bigsqcup_{i=1}^d \sigma_i$, and by recursively using Lemma B.2 we get that

$$\begin{aligned} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma)} &= [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)} \circ \dots \\ &\quad \dots \circ [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_d)} \end{aligned}$$

It then suffices to show the equality in distribution for just one swap operation, so that we can recursively apply the swapped identity mappings while keeping the equality in distribution. We then need to prove that :

$$[\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)}(\mathbf{X}, \tilde{\mathbf{X}}) \stackrel{d}{=} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}](\mathbf{X}, \tilde{\mathbf{X}})$$

Equivalently, if we condition on $\mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1}$ we need to show that

$$\begin{aligned} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)}(\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \\ \stackrel{d}{=} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}](\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \end{aligned}$$

Here crucially we use the fact that σ_1 is a local swap. Indeed, whenever we condition on $\mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1}$, the input values to the mapping $[\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)}$ can be seen as constant with respect to $\mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1}$. Given that σ_1 can only be equal to \emptyset or $\{1\}$, and that therefore its value is determined by

\mathbf{X}_{-1} , hence constant when we condition on \mathbf{X}_{-1} , we get that either

$$\begin{aligned} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)}(\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \\ = [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}](\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \end{aligned}$$

which therefore holds also in distribution or

$$\begin{aligned} [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma_1)}(\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \\ = [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\{1\})}(\mathbf{X}, \tilde{\mathbf{X}}) | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \end{aligned}$$

In that case, it simplifies into

$$X_1, \tilde{X}_1 | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \stackrel{d}{=} \tilde{X}_1, X_1 | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1}$$

which is a consequence of the fact that $\mathbf{X}, \tilde{\mathbf{X}}$ satisfy exchangeability, hence the result.

To prove $[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y \stackrel{d}{=} [\mathbf{X}, \tilde{\mathbf{X}}], Y$, we assume that for every $\mathbf{x} \in \mathbb{R}^d$, $\sigma(\mathbf{x}) \subset \mathcal{H}_0^0(\mathbf{x})$. Jointly taking Y with $(\mathbf{X}, \tilde{\mathbf{X}})$, the proof is the same up to proving that the following holds whenever $1 \in \mathcal{H}_0^0(\mathbf{X})$:

$$X_1, \tilde{X}_1, Y | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1} \stackrel{d}{=} \tilde{X}_1, X_1, Y | \mathbf{X}_{-1}, \tilde{\mathbf{X}}_{-1}$$

By the properties of \mathcal{H}_0^0 , $1 \in \mathcal{H}_0^0(\mathbf{X})$ holds regardless of the value of X_1 when conditioning on \mathbf{X}_{-1} . Now if we write down the densities (as we assumed that the joint distribution has a positive density with respect to a product measure):

$$\begin{aligned} p(\mathbf{x}, \tilde{\mathbf{x}}, y) &= p(y | \mathbf{x}, \tilde{\mathbf{x}}) p(\mathbf{x}, \tilde{\mathbf{x}}) \\ &= p(y | \mathbf{x}) p(\mathbf{x}, \tilde{\mathbf{x}}) \quad \text{as } \tilde{\mathbf{X}} \perp\!\!\!\perp Y | \mathbf{X} \\ &= p(y | \mathbf{x}_{-1}) p(\mathbf{x}, \tilde{\mathbf{x}}) \quad \text{as } X_1 \perp\!\!\!\perp Y | \mathbf{X}_{-1} = \mathbf{x}_{-1} \\ &= p(y | \mathbf{x}_{-1}) p([\mathbf{x}, \tilde{\mathbf{x}}]_{\text{swap}(\{1\})}) \quad \text{by exchangeability} \\ &= p([\mathbf{x}, \tilde{\mathbf{x}}]_{\text{swap}(\{1\})}, y) \end{aligned}$$

Hence the result. \square

B.2. Proof of Proposition 5.2

Proof. Fix σ local swap, $r > 0$, and $\mathbf{z} \in A^r$. Let $S := \sigma(\mathbf{z})$, by definition of local importance scores we have that

$$[(\mathbf{T}_\mu(\cdot), \tilde{\mathbf{T}}_\mu(\cdot))]_{\text{swap}(S)} = [(\mathbf{T}_{\mu_{\text{swap}(S)}}(\cdot), \tilde{\mathbf{T}}_{\mu_{\text{swap}(S)}}(\cdot))]_{\text{swap}(S)}$$

By definition of r -local importance scores, the value of $[\mathbf{T}_{\mu_{\text{swap}(\sigma)}(\mathbf{z})}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(\sigma)}(\mathbf{z})}]$ depends on $\mu_{\text{swap}(\sigma)}$ only through $B(\mathbf{z}, r) \times B(\mathbf{z}, r) \times \mathbb{R}$. Therefore if $\mu_{\text{swap}(\sigma)}$ and $\mu_{\text{swap}(S)}$ coincide on $B(\mathbf{z}, r) \times B(\mathbf{z}, r) \times \mathbb{R}$, then we have

$$\begin{aligned} [\mathbf{T}_{\mu_{\text{swap}(\sigma)}(\mathbf{z})}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(\sigma)}(\mathbf{z})}] \\ = [\mathbf{T}_{\mu_{\text{swap}(S)}(\mathbf{z})}, \tilde{\mathbf{T}}_{\mu_{\text{swap}(S)}(\mathbf{z})}] \\ = [(\mathbf{T}_\mu(\mathbf{z}), \tilde{\mathbf{T}}_\mu(\mathbf{z}))]_{\text{swap}(\sigma(\mathbf{z}))} \end{aligned}$$

Given that $\mathbf{z} \in A^r$, we have that $\forall \mathbf{u}, \mathbf{v} \in B(\mathbf{z}, r)$, $\sigma(\mathbf{u}) = S$. We want to show that, for any $y \in \mathbb{R}$,

$$\begin{aligned}
 & \mu_{\text{swap}(\sigma)}(\mathbf{u}, \mathbf{v}, y) = \mu_{\text{swap}(S)}(\mathbf{u}, \mathbf{v}, y) \\
 \Leftrightarrow & \mathbb{P}([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y = \mathbf{u}, \mathbf{v}, y) \\
 & = \mathbb{P}([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, Y = \mathbf{u}, \mathbf{v}, y) \\
 \Leftrightarrow & \mathbb{P}([\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma)}(\mathbf{X}, \tilde{\mathbf{X}}), Y = \mathbf{u}, \mathbf{v}, y) \\
 & = \mathbb{P}([\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(S)}(\mathbf{X}, \tilde{\mathbf{X}}), Y = \mathbf{u}, \mathbf{v}, y) \\
 \Leftrightarrow & \mathbb{P}(\mathbf{X}, \tilde{\mathbf{X}}, Y = [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(\sigma)}(\mathbf{u}, \mathbf{v}), y) \\
 & = \mathbb{P}(\mathbf{X}, \tilde{\mathbf{X}}, Y = [\mathbf{F}^{Id}, \tilde{\mathbf{F}}^{Id}]_{\text{swap}(S)}(\mathbf{u}, \mathbf{v}), y) \\
 \Leftrightarrow & \mathbb{P}(\mathbf{X}, \tilde{\mathbf{X}}, Y = [\mathbf{u}, \mathbf{v}]_{\text{swap}(\sigma(\mathbf{u}))}, y) \\
 & = \mathbb{P}(\mathbf{X}, \tilde{\mathbf{X}}, Y = [\mathbf{u}, \mathbf{v}]_{\text{swap}(S)}, y)
 \end{aligned}$$

Hence the result. \square

B.3. Extended Flip-Sign Property for Local Swaps and Local Importance Scores

Proposition B.3. σ^r is a local swap and if $\sigma \subset \mathcal{H}_0^0$, then $\sigma^r \subset \mathcal{H}_0^0$. Furthermore,

$$\begin{aligned}
 & [\mathbf{T}_{\hat{\mu}_n}(\mathbf{z}_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(\mathbf{z}_l)]_{1 \leq l \leq L} \\
 & \stackrel{d}{=} [[\mathbf{T}_{\hat{\mu}_n}(\mathbf{z}_l), \tilde{\mathbf{T}}_{\hat{\mu}_n}(\mathbf{z}_l)]_{\text{swap}(\sigma^r(\mathbf{z}_l))}]_{1 \leq l \leq L}
 \end{aligned}$$

Proof. Let \mathbf{x}, \mathbf{z} such that $\mathbf{x}_{[d] \setminus \sigma^r(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma^r(\mathbf{x})}$. We want to show that $\sigma^r(\mathbf{x}) = \sigma^r(\mathbf{z})$.

If $\mathbf{x} \notin A^r$, then $\sigma^r(\mathbf{x}) = \emptyset$, so $\mathbf{x} = \mathbf{z}$ and $\sigma^r(\mathbf{x}) = \sigma^r(\mathbf{z})$.

If $\mathbf{x} \in A^r$, let us first show that it implies $\mathbf{z} \in A^r$. Let $\mathbf{y} \in B(\mathbf{z}, r)$, show that $\sigma(\mathbf{y}) = \sigma(\mathbf{z})$. We have that $\mathbf{y} - (\mathbf{z} - \mathbf{x}) \in B(\mathbf{x}, r)$, and given that $\mathbf{x} \in A^r$, we get $\sigma(\mathbf{x}) = \sigma(\mathbf{y} - (\mathbf{z} - \mathbf{x}))$. As σ is a local swap, we have $\sigma(\mathbf{x}) = \sigma(\mathbf{z})$, and we also get $\sigma(\mathbf{y} - (\mathbf{z} - \mathbf{x})) = \sigma(\mathbf{y})$ because

$$\begin{aligned}
 & [\mathbf{y} - (\mathbf{z} - \mathbf{x})]_{[d] \setminus \sigma(\mathbf{y} - (\mathbf{z} - \mathbf{x}))} = \mathbf{y}_{[d] \setminus \sigma(\mathbf{y} - (\mathbf{z} - \mathbf{x}))} \\
 \Leftrightarrow & [\mathbf{y} - (\mathbf{z} - \mathbf{x})]_{[d] \setminus \sigma(\mathbf{x})} = \mathbf{y}_{[d] \setminus \sigma(\mathbf{x})} \\
 \Leftrightarrow & [(\mathbf{z} - \mathbf{x})]_{[d] \setminus \sigma(\mathbf{x})} = 0 \\
 \Leftrightarrow & \mathbf{x}_{[d] \setminus \sigma^r(\mathbf{x})} = \mathbf{z}_{[d] \setminus \sigma^r(\mathbf{x})}
 \end{aligned}$$

We can now conclude: $\sigma^r(\mathbf{x}) = \sigma(\mathbf{x})$ as $\mathbf{x} \in A^r$, and given that σ is a local swap we get that $\sigma(\mathbf{x}) = \sigma(\mathbf{z})$. Finally, as $\mathbf{z} \in A^r$, we have $\sigma(\mathbf{z}) = \sigma^r(\mathbf{z})$ and therefore $\sigma^r(\mathbf{x}) = \sigma^r(\mathbf{z})$.

Assume $\sigma \subset \mathcal{H}_0^0$. Fix $\mathbf{z} \in \mathbb{R}^d$, assume that $\sigma^r(\mathbf{z}) \neq \emptyset$ and take $j \in \sigma^r(\mathbf{z})$. That implies $\mathbf{z} \in A^r$, and therefore $\forall \mathbf{y} \in B(\mathbf{z}, r)$ we have $j \in \sigma^r(\mathbf{z}) = \sigma(\mathbf{z}) = \sigma(\mathbf{y}) \subset \mathcal{H}_0^0(\mathbf{y})$. Therefore $j \in \bigcap_{\mathbf{y} \in B(\mathbf{z}, r)} \mathcal{H}_0^0(\mathbf{y}) = \mathcal{H}_0^r(\mathbf{z})$. We then conclude that $\sigma^r \subset \mathcal{H}_0^r$.

The last statement is a concatenation of Proposition 5.2 and the fact that $\mu = \mu_{\text{swap}(\sigma^r)}$. \square

C. Semi-synthetic Data Experiments

Our simulations previously described were entirely based on synthetic data. Alternatively, using real SNPs data and then fitting a HMM model yields the same experimental results, which we did for data from the 1000 Genomes Project (Consortium et al., 2015), where we obtained around 2000 individual samples for 27 distinct segments of chromosome 19 containing an average of 50 SNPs per segment, and filtered out SNPs that are extremely correlated (above 0.95). This is because HMMs can truly capture the covariate distribution of a SNP dataset and are a good model for a downstream feature selection with the knockoff procedure. For simplicity, and in order to scale with the number of samples (which is limited with real data), we described simulations based on synthetic covariates.

D. Saliency-based Partitioning

Saliency maps (Lipton, 2016) have emerged as a popular tool for interpretability in Neural Networks. A saliency map allows to identify, under a trained model that minimizes some loss L , which input variation has the strongest impact on the loss at a given training point. Training a neural network on the concatenated vector $\mathbf{X}, \tilde{\mathbf{X}}$ to predict Y , saliency maps can be used as importance scores at any given training point. That is, denoting $g_\theta : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ a classifier parametrized by θ (such as a neural network), consider $\hat{\theta}_n$ the output of training such model on the actual data. We can now compute the saliency scores:

$$\mathbb{S}, \tilde{\mathbb{S}} : \begin{cases} \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d \\ \mathbf{X}_i, \tilde{\mathbf{X}}_i \mapsto \nabla_{\mathbf{X}, \tilde{\mathbf{X}}} L(g_{\hat{\theta}_n}(\mathbf{X}_i, \tilde{\mathbf{X}}_i) - Y_i) \end{cases}$$

Notice that the saliency scores are only computed for training points $\mathbf{X}_i, \tilde{\mathbf{X}}_i, Y_i$, but our definition of $\mathbb{S}, \tilde{\mathbb{S}}$ assumes that we can expand the saliency scores to the whole feature space (for example, through smoothing). The reason why these mappings can not be immediately used as local importance scores is because of the training process: the output of the training process are the parameter estimates in $\hat{\theta}_n = \hat{\theta}_n(\mathbf{X}, \tilde{\mathbf{X}}, Y)$, which are constructed based on the global training set. Even though the saliency is local at a point the swap operation can not go through the training process, i.e. we can not relate $g_{\hat{\theta}_n(\mathbf{X}, \tilde{\mathbf{X}}, Y)}([\mathbf{x}, \tilde{\mathbf{x}}]_{\text{swap}(\sigma)})$ and $g_{\hat{\theta}_n([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(\sigma)}, Y)}([\mathbf{x}, \tilde{\mathbf{x}}])$. This is due to the influence that a training point lying in one region of the space may have at another point (during evaluation) at a different region (Koh & Liang, 2017). Still, we can use these saliency scores to partition the space based on a subsample of the whole initial dataset, and then run the Knockoff procedure with local importance scores at points located in each subregion. This method has the advantage of being computationally less expensive than the previous one, especially in high dimensions.