

A Theoretical details

A.1 A note about the assumptions

Note about the assumptions In [theorem 1](#), assumption 1 consists of three parts that can all be validated on observed data: 1) that the gradient flow converges, 2) that the confounder value of the surrogate matches the confounder value whose effect is of interest, and 3) that the surrogate intervention lies in the support of the pre-outcome variables. Assumption 2 is required for expectations and their gradients to exist and be finite. In [theorem 2](#), assumption 1 requires a consistent estimator of $\mathbb{E}[\mathbf{y} | \mathbf{t}]$, which can be provided with regression. Assumption 3 lists regularity conditions which help control how the surrogate estimation error propagates to the effect error.

A.2 Proof of [Theorem 1](#)

We restate the theorem for completeness:

Theorem 1. *Assume C-REDUNDANCY holds. Assuming the following:*

1. Let $\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*))$ be the limiting solution to the gradient flow equation $\frac{d\tilde{\mathbf{t}}(s)}{ds} = -\nabla_{\tilde{\mathbf{t}}}(\mathbf{h}(\tilde{\mathbf{t}}(s)) - \mathbf{h}(\mathbf{t}_2^*))^2$, initialized at $\tilde{\mathbf{t}}(0) = \mathbf{t}^*$; i.e. $\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)) = \lim_{s \rightarrow \infty} \tilde{\mathbf{t}}(s)$. Further, let $\mathbf{h}(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*))) = \mathbf{h}(\mathbf{t}_2^*)$ and $\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)) \in \text{supp}(\mathbf{t})$.
2. $f(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}), \boldsymbol{\eta})$ and $\mathbf{h}(\tilde{\mathbf{t}})$ as functions of $\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}})$ are continuous and differentiable and the derivatives exist for all $\tilde{\mathbf{t}}, \boldsymbol{\eta}$. Let $\nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}), \boldsymbol{\eta})$ exist and be bounded and integrable w.r.t. the probability measure corresponding to $\mathbf{p}(\boldsymbol{\eta})$, for all values of $\tilde{\mathbf{t}}$ and $h(\tilde{\mathbf{t}})$.

Then the conditional effect (and therefore the average effect) is identified:

$$\phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) = \phi(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)), \mathbf{h}(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)))) = \mathbb{E}[\mathbf{y} | \mathbf{t} = \mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*))] \quad (10)$$

Proof. Recall definition of conditional effect $\phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2)) = \mathbb{E}_{\boldsymbol{\eta}} f(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2), \boldsymbol{\eta})$. Recall $\nabla_{\tilde{\mathbf{t}}}$ is the gradient with respect to the first argument of f , that is $\tilde{\mathbf{t}}$. First, by assumption 2, \mathbb{E} and ∇ commute, under the dominated convergence theorem. Then, by C-REDUNDANCY

$$\nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*))^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = \nabla_{\tilde{\mathbf{t}}} \mathbb{E}_{\boldsymbol{\eta}} f(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*), \boldsymbol{\eta})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = \mathbb{E}_{\boldsymbol{\eta}} [\nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*), \boldsymbol{\eta})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}})] = 0.$$

Now consider the gradient flow equation $d\tilde{\mathbf{t}}(s)/ds = -\nabla_{\tilde{\mathbf{t}}}(\mathbf{h}(\tilde{\mathbf{t}}) - \mathbf{h}(\mathbf{t}_2^*))^2$. We refer to the gradient evaluated at $\tilde{\mathbf{t}}$ as $\Delta\tilde{\mathbf{t}} = -\nabla_{\tilde{\mathbf{t}}}(\mathbf{h}(\tilde{\mathbf{t}}) - \mathbf{h}(\mathbf{t}_2^*))^2 = -2(\mathbf{h}(\tilde{\mathbf{t}}) - \mathbf{h}(\mathbf{t}_2^*))\nabla_{\tilde{\mathbf{t}}}\mathbf{h}(\tilde{\mathbf{t}})$. We will express $\phi(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)), \mathbf{h}(\mathbf{t}_2^*))$ as defined by the starting point $\phi(\mathbf{t}^*, h(\mathbf{t}_2^*))$ and the gradient flow equation.

Let the solution path to the gradient flow equation be C with $\mathbf{t}^*, \mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*))$ being the starting and ending points respectively. By the Gradient Theorem [\[26\]](#), we have that $\phi(\mathbf{t}^*, h(\mathbf{t}_2^*))$ and $\phi(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)), \mathbf{h}(\mathbf{t}_2^*))$ are related via the line integral over C :

$$\int_C \nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*)) \cdot d\tilde{\mathbf{t}} = \phi(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)), \mathbf{h}(\mathbf{t}_2^*)) - \phi(\mathbf{t}^*, h(\mathbf{t}_2^*))$$

Let $\tilde{\mathbf{t}}(s)$ be a parametrization of solution path C by the scalar time $s \in [0, \infty)$. Now, to obtain the value of $\phi(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*))$, we will compute the line integral over the vector field defined by $\nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*))$, which exists by assumption 2 in [theorem 1](#), evaluated along the path C defined by $\Delta\tilde{\mathbf{t}}(s)$:

$$\begin{aligned} \phi(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)), \mathbf{h}(\mathbf{t}_2^*)) &= \phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) + \int_C \nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}, h(\mathbf{t}_2^*)) \cdot d\tilde{\mathbf{t}} \\ &= \phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) + \int_0^\infty \nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}(s), h(\mathbf{t}_2^*))^\top \frac{d\tilde{\mathbf{t}}(s)}{ds} ds \\ &= \phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) + \int_0^\infty \nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}(s), h(\mathbf{t}_2^*))^\top \Delta\tilde{\mathbf{t}}(s) ds \\ &= \phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) \\ &\quad + \int_0^\infty -2((\mathbf{h}(\tilde{\mathbf{t}}(s)) - \mathbf{h}(\mathbf{t}_2^*))) \nabla_{\tilde{\mathbf{t}}} \phi(\tilde{\mathbf{t}}(s), h(\mathbf{t}_2^*))^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}(s)) ds \\ &= \phi(\mathbf{t}^*, h(\mathbf{t}_2^*)) + 0 \quad \{\text{by C-REDUNDANCY}\} \end{aligned} \quad (11)$$

Finally, by assumption 1 in [theorem 1](#), $h(\tau'(\tau^*, h(\tau_2^*))) = h(\tau_2^*)$, and so

$$\phi(\tau^*, h(\tau_2^*)) = \phi(\tau'(\tau^*, h(\tau_2^*)), h(\tau_2^*)) = \phi(\tau'(\tau^*, h(\tau_2^*)), h(\tau'(\tau^*, h(\tau_2^*)))) \quad (12)$$

For clarity, the same equation, but using τ' and suppressing dependence on $\tau^*, h(\tau_2^*)$:

$$\phi(\tau^*, h(\tau_2^*)) = \phi(\tau', h(\tau_2^*)) = \phi(\tau', h(\tau')) \quad (13)$$

Under the causal model for EFC, the outcome $\mathbf{y} = f(\mathbf{t}, h(\mathbf{t}), \boldsymbol{\eta})$. Then, $\forall \tilde{\mathbf{t}} \in \text{supp}(\mathbf{p}(\mathbf{t}))$,

$$\mathbb{E}[\mathbf{y} | \mathbf{t} = \tilde{\mathbf{t}}] = \mathbb{E}_{\boldsymbol{\eta}}[f(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}), \boldsymbol{\eta})] = \phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}})). \quad (14)$$

Using that $\tau'(\tau^*, \tau_2^*) \in \text{supp}(\mathbf{p}(\mathbf{t}))$ and [eqs. \(13\) and \(14\)](#), the conditional effect is identified

$$\begin{aligned} \phi(\tau^*, h(\tau_2^*)) &= \phi(\tau'(\tau^*, h(\tau_2^*)), h(\tau'(\tau^*, h(\tau_2^*)))) \\ &= \mathbb{E}[\mathbf{y} | \mathbf{t} = \tau'(\tau^*, h(\tau_2^*))] \end{aligned} \quad (15)$$

Thus, the conditional effect, and consequently the average effect, are identified as $\mathbb{E}[\mathbf{y} | \tau'(\tau^*, h(\tau_2^*))]$ and $\tau(\tau^*) = \mathbb{E}_{h(\mathbf{t})} \mathbb{E}[\mathbf{y} | \tau'(\tau^*, h(\mathbf{t}))]$ respectively. \square

Note about convergence of gradient flow Any ODE's solution, if it exists and converges, converges to an ω -limit set [\[27\]](#). An ω -limit set is nonempty when the solution path lies entirely in a closed and bounded set and can consist of limit cycles, equilibrium points, or neither [\[13, 27\]](#). A gradient flow equation $d\tilde{\mathbf{t}}(s)/ds = -\nabla h(\tilde{\mathbf{t}})$ (also called a gradient system) has the special property that its ω -limit set only consists of critical points of $h(\tilde{\mathbf{t}})$; critical points of $h(\tilde{\mathbf{t}})$ are also equilibrium points of the gradient flow equation [\[13\]](#). Further, if $\nabla h(\tilde{\mathbf{t}})$ exists and is bounded and $h(\tilde{\mathbf{t}})$ has bounded sublevel sets ($\{\tilde{\mathbf{t}} : h(\tilde{\mathbf{t}}) \leq c\}$), then the solution to the gradient flow equation will entirely lie within a bounded set. This is because along the solution path, $h(\tilde{\mathbf{t}}(s))$ always decreases meaning that the solution will remain in any sublevel set it started in. Thus, if $h(\tilde{\mathbf{t}})$ has bounded sublevel sets, the solution of the gradient flow equation will converge only to critical points of $h(\tilde{\mathbf{t}})$.

A.3 Estimation error in LODE

Theorem 2. Consider the conditional effect $\phi(\tau^*, h(\tau_2^*))$. Let $\hat{\mathbf{t}}(\tau^*, h(\tau_2^*))$ be the estimate of the surrogate intervention computed by LODE, computed via Euler integration of the gradient flow $\frac{d\tilde{\mathbf{t}}(s)}{ds} = -\nabla_{\tilde{\mathbf{t}}} (h(\tilde{\mathbf{t}}(s)) - h(\tau_2^*))^2$, initialized at $\tilde{\mathbf{t}}(0) = \tau^*$. Assume the true surrogate $\tau'(\tau^*, h(\tau_2^*))$ exists and is the limiting solution to the gradient flow equation.

1. Let the finite sample estimator of $\mathbb{E}[\mathbf{y} | \mathbf{t} = \tilde{\mathbf{t}}]$ be $\hat{f}(\tilde{\mathbf{t}})$. Let the error for all $\tilde{\mathbf{t}}$ be bounded, $|\hat{f}(\tilde{\mathbf{t}}) - \mathbb{E}[\mathbf{y} | \mathbf{t} = \tilde{\mathbf{t}}]| \leq c(N)$, where N is the sample size and $\lim_{N \rightarrow \infty} c(N) = 0$.
2. Assume K Euler integrator steps were taken to find the surrogate estimate $\hat{\mathbf{t}}(\tau^*, h(\tau_2^*))$, each of size ℓ . Let the maximum confounder mismatch be $\max_{i \leq K} (h(\tilde{\mathbf{t}}_i) - h(\tau_2^*))^2 = M$.
3. Let $L_{z, \tilde{\mathbf{t}}}$ be the Lipschitz-constant of $\phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2))$ as a function of $h(\tilde{\mathbf{t}}_2)$, for fixed $\tilde{\mathbf{t}}$. Let L_e be the Lipschitz-constant of $\mathbb{E}[\mathbf{y} | \mathbf{t} = \tilde{\mathbf{t}}] = \phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}))$ as a function of $\tilde{\mathbf{t}}$. Assume h has a gradient with bounded norm, $\|\nabla h(\tilde{\mathbf{t}})\|_2 \leq L_h$. Assume f 's Hessian has bounded eigenvalues: $\forall \tilde{\mathbf{t}}, \tilde{\mathbf{t}}_2, \|\nabla_{\tilde{\mathbf{t}}}^2 \phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2))\|_2 \leq \sigma_{H\phi}$.

The conditional effect estimate error, $\xi(\tau^*, h(\tau_2^*)) = |\hat{f}(\hat{\mathbf{t}}) - \phi(\tau^*, h(\tau_2^*))|$, is upper bounded by:

$$c(N) + \min (L_e \|\tau' - \hat{\mathbf{t}}\|_2, 2K\ell^2 (\mathcal{O}(\ell) + M\sigma_{H\phi} L_h^2) + L_{z, \hat{\mathbf{t}}} \|h(\hat{\mathbf{t}}) - h(\tau_2^*)\|_2) \quad (16)$$

Proof. (of [Theorem 2](#)) Recall the definition of conditional effect: $\phi(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2)) = \mathbb{E}_{\boldsymbol{\eta}} f(\tilde{\mathbf{t}}, h(\tilde{\mathbf{t}}_2), \boldsymbol{\eta})$.

LODE's estimate of the conditional effect is $\hat{f}(\hat{\mathbf{t}}(\tau^*, h(\tau_2^*)))$. We will suppress notation for dependence on $\tau^*, h(\tau_2^*)$, and use τ' and $\hat{\mathbf{t}}$ to refer to the true surrogate intervention and the estimated surrogate interventions respectively. Note \hat{f} is the estimate of the conditional expectation $\mathbb{E}[\mathbf{y} | \mathbf{t} = \tilde{\mathbf{t}}]$, learned from N samples. We first bound the error by splitting into two parts and bounding each separately:

$$\begin{aligned} |\xi(\tau^*, h(\tau_2^*))| &= |\hat{f}(\hat{\mathbf{t}}) - \phi(\tau^*, h(\tau_2^*))| \\ &\leq |\hat{f}(\hat{\mathbf{t}}) - \phi(\hat{\mathbf{t}}, h(\hat{\mathbf{t}}))| + |\phi(\hat{\mathbf{t}}, h(\hat{\mathbf{t}})) - \phi(\tau^*, h(\tau_2^*))| \\ &\leq c(N) + |\phi(\hat{\mathbf{t}}, h(\hat{\mathbf{t}})) - \phi(\tau^*, h(\tau_2^*))| \\ &\leq |\phi(\hat{\mathbf{t}}, h(\hat{\mathbf{t}})) - \phi(\hat{\mathbf{t}}, h(\tau_2^*))| + |\phi(\hat{\mathbf{t}}, h(\tau_2^*)) - \phi(\tau^*, h(\tau_2^*))| + c(N) \end{aligned}$$

The first term is bounded via the Lipschitz-ness of ϕ as a function of $h(\tilde{\tau})$ with fixed first argument $\tilde{\tau} = \hat{\tau}$.

$$|\phi(\hat{\tau}, h(\hat{\tau})) - \phi(\hat{\tau}, h(\tau_2^*))| \leq L_{z,\hat{\tau}} \|h(\hat{\tau}) - h(\tau_2^*)\|$$

We now bound the remaining term. Recall that LODE's computation of the surrogate intervention involved K gradient steps, each of size ℓ . We work with a constant step-size but the analysis can be generalized to a non-uniform step size. Indexing steps with i , let $d_i = h(\tilde{\tau}_i) - h(\tau_2^*)$ be the confounder mismatch error at the i th iterate. Then note that $\hat{\tau} = \tau^* - \ell \sum_{i=0}^{K-1} 2d_i \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)$. We can use this to bound the error $\phi(\hat{\tau}, h(\tau_2^*)) - \phi(\tau^*, h(\tau_2^*))$. With $\tilde{\tau}_K = \hat{\tau}$ and $\tilde{\tau}_0 = \tau^*$, we proceed by expressing the error as a telescoping sum and using the Taylor expansion for $\phi(\tilde{\tau}, h(\tau_2^*))$ in terms of the the first argument $\tilde{\tau}$.

$$\phi(\hat{\tau}, h(\tau_2^*)) - \phi(\tau^*, h(\tau_2^*)) = \sum_{i=0}^{K-1} \phi(\tilde{\tau}_{i+1}, h(\tau_2^*)) - \phi(\tilde{\tau}_i, h(\tau_2^*)) \quad (17)$$

$$= \sum_{i=0}^{K-1} \nabla_{\tilde{\tau}} \phi(\tilde{\tau}_i, h(\tau_2^*))^\top (\tilde{\tau}_{i+1} - \tilde{\tau}_i) \quad (18)$$

$$+ \frac{1}{2} (\tilde{\tau}_{i+1} - \tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) (\tilde{\tau}_{i+1} - \tilde{\tau}_i) + \mathcal{O}(\|\tilde{\tau}_{i+1} - \tilde{\tau}_i\|_2^3) \quad (19)$$

$$= \sum_{i=0}^{K-1} 2\ell d_i \nabla_{\tilde{\tau}} \phi(\tilde{\tau}_i, h(\tau_2^*))^\top \nabla_{\tilde{\tau}} h(\tilde{\tau}_i) + 2(\ell d_i)^2 \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) \nabla_{\tilde{\tau}} h(\tilde{\tau}_i) + \mathcal{O}(\ell^3) \quad (20)$$

$$= \sum_{i=0}^{K-1} 0 + 2(\ell d_i)^2 \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) \nabla_{\tilde{\tau}} h(\tilde{\tau}_i) + \mathcal{O}(\ell^3) \quad (21)$$

$$= \mathcal{O}(K\ell^3) + \sum_{i=0}^{K-1} 2(\ell d_i)^2 \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) \nabla_{\tilde{\tau}} h(\tilde{\tau}_i) \quad (22)$$

$$\leq \mathcal{O}(K\ell^3) + \sum_{i=0}^{K-1} 2(\ell(h(\tilde{\tau}_i) - h(\tau_2^*)))^2 |\nabla_{\tilde{\tau}} h(\tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)| \quad (23)$$

$$\leq \mathcal{O}(K\ell^3) + \sum_{i=0}^{K-1} 2\ell^2 M |\nabla_{\tilde{\tau}} h(\tilde{\tau}_i)^\top \nabla_{\tilde{\tau}}^2 \phi(\tilde{\tau}_i, h(\tau_2^*)) \nabla_{\tilde{\tau}} h(\tilde{\tau}_i)| \quad (24)$$

$$\leq \mathcal{O}(K\ell^3) + \sum_{i=0}^{K-1} 2\ell^2 M \sigma_{\mathbb{H}\phi} \|\nabla_{\tilde{\tau}} h(\tilde{\tau}_i)\|_2^2 \quad (25)$$

$$\leq \mathcal{O}(K\ell^3) + \sum_{i=0}^{K-1} 2\ell^2 M \sigma_{\mathbb{H}\phi} L_h^2 \quad (26)$$

$$= 2K\ell^2 (\mathcal{O}(\ell) + M\sigma_{\mathbb{H}\phi} L_h^2), \quad (27)$$

where the inequalities follow by the maximum value of $(h(\tilde{\tau}_i) - h(\tau_2^*))^2$, bounded eigenvalues of the Hessian of ϕ and the Lipschitz-ness of $h(\tilde{\tau})$.

Another way we bound the error is via the Lipschitz constant of the conditional expectation as a function of $\tilde{\tau}$. Recall this is L_e . An alternate bound on the error is as follows:

$$|\phi(\hat{\tau}, h(\hat{\tau})) - \phi(\tau^*, h(\tau_2^*))| = |\phi(\hat{\tau}, h(\hat{\tau})) - \phi(\tau', h(\tau'))| \leq L_e \|\tau' - \hat{\tau}\|_2$$

The bound follows:

$$|\xi(\tilde{\tau}, h(\tau_2^*))| \leq c(N) + \min(L_e \|\tau' - \hat{\tau}\|_2, \quad 2K\ell^2 (\mathcal{O}(\ell) + M\sigma_{\mathbb{H}\phi} L_h^2) + L_{z,\hat{\tau}} \|h(\hat{\tau}) - h(\tau_2^*)\|_2)$$

□

A.3.1 A note on linear confounder functions and LODE

In the proof above, the error in Euler integration accumulates due to terms like this one: $\nabla_{\tilde{\mathbf{t}}}^\top \mathbf{h}(\tilde{\mathbf{t}}) \nabla_{\tilde{\mathbf{t}}}^2 f(\tilde{\mathbf{t}}, \mathbf{h}(\mathbf{t}^*), \boldsymbol{\eta}) \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}})$. For a linear confounder function that satisfies $\nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = \boldsymbol{\beta}$, such terms can be expressed as $\boldsymbol{\beta}^\top \nabla_{\tilde{\mathbf{t}}} (\nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, \mathbf{h}(\mathbf{t}^*), \boldsymbol{\eta})^\top \boldsymbol{\beta}) = \boldsymbol{\beta}^\top \nabla_{\tilde{\mathbf{t}}} (0) = 0$ under C-REDUNDANCY. Thus, such error does not accumulate even with large step sizes.

Further, note that the gradient flow equation in LODE for the causal model A in section 4 is a linear ODE whose solution has a closed form expression and one can estimate the surrogate without numerical integration [27].

A.4 Proof of sufficiency of Effect Connectivity

Theorem 3. *Under Effect Connectivity, eq. (9), any surrogate intervention $\mathbf{t}'(\mathbf{t}^*, \mathbf{h}(\mathbf{t}_2^*)) \in \text{supp}(\mathbf{t})$.*

Proof. Recall $\phi(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}})) = \mathbb{E}_{\boldsymbol{\eta}} f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}), \boldsymbol{\eta})$. We have $\forall \mathbf{t}^* \in \text{supp}(\mathbf{p}(\mathbf{t}))$:

$$\mathbf{p}(\mathbf{h}(\mathbf{t}) = \mathbf{h}(\mathbf{t}_2^*)) > 0 \implies \mathbf{p}(\phi(\mathbf{t}, \mathbf{h}(\mathbf{t})) = \phi(\mathbf{t}^*, \mathbf{h}(\mathbf{t}_2^*)) \mid \mathbf{h}(\mathbf{t}) = \mathbf{h}(\mathbf{t}_2^*)) > 0.$$

This implies $\exists \mathbf{t}' \in \text{supp}(\mathbf{t})$, $\phi(\mathbf{t}', \mathbf{h}(\mathbf{t}_2^*)) = \phi(\mathbf{t}^*, \mathbf{h}(\mathbf{t}_2^*))$, s.t. $\mathbf{h}(\mathbf{t}') = \mathbf{h}(\mathbf{t}_2^*)$.

Then, $\phi(\mathbf{t}^*, \mathbf{h}(\mathbf{t}_2^*)) = \phi(\mathbf{t}', \mathbf{h}(\mathbf{t}_2^*)) = \phi(\mathbf{t}', \mathbf{h}(\mathbf{t}')) = \mathbb{E}[\mathbf{y} \mid \mathbf{t} = \mathbf{t}']$. \square

A.5 Necessity of Effect Connectivity for Nonparametric effect estimation in EFC

Theorem 4. *Effect Connectivity is necessary for nonparametric effect estimation in EFC.*

Proof. (Proof of Theorem 4) Let the outcome be $\mathbf{y} = f(\mathbf{t}, \mathbf{h}(\mathbf{t}))$. Recall the joint distribution $\mathbf{p}(\mathbf{t}, \mathbf{y})$ and let $\mathbf{h}(\mathbf{t})$ be the confounder. Let Effect Connectivity be violated, i.e. there exists a non-measure-zero subset $B \in \text{supp}(\mathbf{t}) \times \text{supp}(\mathbf{h}(\mathbf{t}))$ such that ⁶:

$$\forall \tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2) \in B, \quad \mathbf{p}(f(\mathbf{t}, \mathbf{h}(\mathbf{t})) = f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) \mid \mathbf{h}(\mathbf{t}) = \mathbf{h}(\tilde{\mathbf{t}}_2)) = 0.$$

Now, we construct a new outcome $\mathbf{y}_2 = f_2(\mathbf{t}, \mathbf{h}(\mathbf{t}))$ and show the conditional effects for this new outcome are different from the one defined by f on $\forall(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) \in B$. Let

$$f_2(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) = f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) + 10 * 1((\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) \in B).$$

We have $f_2(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}})) = f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}})) \forall \tilde{\mathbf{t}} \in \text{supp}(\mathbf{t})$, as the additional term in f_2 is only present for $(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) \in B$; this follows from the fact that $\forall \tilde{\mathbf{t}} \in \text{supp}(\mathbf{t})$, $(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}})) \notin B$ as

$$\mathbf{p}[f(\mathbf{t}, \mathbf{h}(\mathbf{t})) = f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}})) \mid \mathbf{h}(\mathbf{t}) = \mathbf{h}(\tilde{\mathbf{t}})] = \mathbf{p}[f(\mathbf{t}, \mathbf{h}(\mathbf{t})) = f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}))] > 0.$$

Thus, $\mathbf{p}(\mathbf{y}, \mathbf{t}) =^d \mathbf{p}(\mathbf{y}_2, \mathbf{t})$ are equal in distribution since $B \cap \text{supp}(\mathbf{t}, \mathbf{h}(\mathbf{t})) = \emptyset$. This means that the conditional effects are different for the outcomes \mathbf{y}, \mathbf{y}_2 for all $(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2)) \in B$:

$$\mathbb{E}[\mathbf{y} \mid \text{do}(\mathbf{t} = \tilde{\mathbf{t}}), \mathbf{h}(\mathbf{t}) = \mathbf{h}(\tilde{\mathbf{t}}_2)] \neq \mathbb{E}[\mathbf{y}_2 \mid \text{do}(\mathbf{t} = \tilde{\mathbf{t}}), \mathbf{h}(\mathbf{t}) = \mathbf{h}(\tilde{\mathbf{t}}_2)]$$

Therefore, for causal models that violates Effect Connectivity, there exist observationally equivalent causal models with different causal effects. Thus, nonparametric effect estimation is impossible. Thus, Effect Connectivity is required for EFC. \square

A.6 Algorithmic details

We give in algorithm 1 pseudocode for LODE.

Extensions of LODE Consider that we have access to $\mathbf{m}(\mathbf{h}(\mathbf{t}))$ for some bijective differentiable function $\mathbf{m}(\cdot)$, instead of $\mathbf{h}(\mathbf{t})$. The orthogonality in C-REDUNDANCY holds $\nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2), \boldsymbol{\eta})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{m}(\mathbf{h}(\tilde{\mathbf{t}})) = \mathbf{m}'(\mathbf{h}(\tilde{\mathbf{t}})) \nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2), \boldsymbol{\eta})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = 0$. Then, using $\mathbf{m}(\mathbf{h}(\tilde{\mathbf{t}}))$ to compute the surrogate $\mathbf{t}'(\mathbf{t}^*, \mathbf{h}(\mathbf{t}_2^*))$, LODE would estimate valid effects. Similarly, LODE can estimate the effect on any differentiable transformation of the outcome $\mathbf{m}(\mathbf{y})$, because $\nabla_{\tilde{\mathbf{t}}} \mathbf{m}(\mathbf{y}_{\tilde{\mathbf{t}}})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = \mathbf{m}'(\mathbf{y}_{\tilde{\mathbf{t}}}) \nabla_{\tilde{\mathbf{t}}} f(\tilde{\mathbf{t}}, \mathbf{h}(\tilde{\mathbf{t}}_2), \boldsymbol{\eta})^\top \nabla_{\tilde{\mathbf{t}}} \mathbf{h}(\tilde{\mathbf{t}}) = 0$ holds.

⁶Non-zero w.r.t. the product measure over $\text{supp}(\mathbf{t}) \times \text{supp}(\mathbf{h}(\mathbf{t}))$ due to \mathbf{p} .

Algorithm 1: LODE for $\text{do}(\mathbf{t} = \mathbf{t}^*)$

Input: Functional confounder $h(\mathbf{t})$; tolerance ϵ

Output: Conditional effects of \mathbf{t}^* , $h(\mathbf{t}_2^*)$

- 1 Regress \mathbf{y} on \mathbf{t} and compute $\hat{f}() := \arg \min_{\mathbf{u} \in \mathcal{F}} \mathbb{E}_{\mathbf{y}, \mathbf{t}} (\mathbf{y} - \mathbf{u}(\mathbf{t}))^2$.
- 2 To estimate effects of \mathbf{t}^* , $h(\mathbf{t}_2^*)$, compute the surrogate intervention $\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*))$ by Euler integrating the gradient flow equation, initialized at $\tilde{\mathbf{t}} = \mathbf{t}^*$, until $(h(\tilde{\mathbf{t}}_s) - h(\mathbf{t}_2^*))^2 < \epsilon$.

$$\frac{d\tilde{\mathbf{t}}(s)}{ds} = \nabla_{\tilde{\mathbf{t}}} (h(\tilde{\mathbf{t}}_s) - h(\mathbf{t}_2^*))^2,$$

- 3 Return $\hat{f}(\mathbf{t}'(\mathbf{t}^*, h(\mathbf{t}_2^*)))$;
-

B Experimental Details

B.1 Functional confounders in GWAS

Here, we show how $h(\mathbf{t}) = \mathbf{A}\mathbf{t}$ and \mathbf{A} reflect the traditional PCA based adjustment in GWAS. Recall population structure acts as a confounder in GWAS. Price et al. [19] demonstrated that using the principal components of the normalized genetic relationships matrix adjusts for confounding due to population structure in GWAS. Let the genotype matrix be \mathbf{G} with people as rows and SNPs as columns, such that each element is one of 0, 1/2, 1, where 1/2 and 1 refer to one and two copies of the allele respectively at the position of the SNP. With p_s as the allele frequency at SNP s [28], Φ is the genetic relationship matrix whose elements are defined as $\Phi_{i,j} = \frac{1}{S} \sum_{s=1}^S (G_{i,s} - p_s)(G_{j,s} - p_s) / p_s(1 - p_s)$. Then, Price et al. [19] compute the top K (10 suggested) principal components of Φ to use as the axes of variation due to the population structure. The eigenvectors of Φ are the left eigenvectors of $\hat{\mathbf{G}}$ such that $\Phi = \hat{\mathbf{G}}\hat{\mathbf{G}}^T$ which capture independent axes of variation of individuals.

Price et al. [19] exploit the idea that if a SNP aligns with some of the axes of variation, this is due to the population structure. These axes of variation are the top K eigenvectors \mathbf{U} of $\phi = \hat{\mathbf{G}}\hat{\mathbf{G}}^T \approx \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, where $\mathbf{U} \in \mathbb{R}^{N \times K}$, $\Phi \in \mathbb{R}^{N \times N}$ and $\mathbf{\Lambda} \in \mathbb{R}^{K \times K}$. Here, \mathbf{U} are also the left singular vectors of $\hat{\mathbf{G}} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ where $\mathbf{\Sigma} \in \mathbb{R}^{K \times K}$ is diagonal, and $\mathbf{V} \in \mathbb{R}^{S \times K}$. We use \approx to denote that the chosen K eigenvectors explain the variation due to population structure; what remains are random mutations.

Let the s th SNP be $\hat{\mathbf{G}}_{:,s} \in \mathbb{R}^N$, which is a column in $\hat{\mathbf{G}}$. In Price et al. [19], population structure in the s th SNP is captured in $\hat{\mathbf{G}}_{:,s}^T \mathbf{U}$. In words, projecting the SNP $\hat{\mathbf{G}}_{:,s}$ onto the axes of variation in individuals gives the population structure between s th SNP and the outcome. This projection $\hat{\mathbf{G}}_{:,s}^T \mathbf{U}$ is a row of $\hat{\mathbf{G}}^T \mathbf{U} \in \mathbb{R}^{S \times K}$. In turn, $\hat{\mathbf{G}}^T \mathbf{U} \in \mathbb{R}^{S \times K}$ is the population structure in all SNPs. Projecting this population structure onto the genotype of an individual gives the confounding due to population structure amongst the SNPs present in the genotype. With $G_{j,\cdot} \in \{0, 1/2, 1\}^S$ as the genotype for an individual j , this projection is $((\hat{\mathbf{G}}^T \mathbf{U})^T G_{j,\cdot})$. However, $\hat{\mathbf{G}} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ implies that $\hat{\mathbf{G}}^T \mathbf{U} \approx \mathbf{V}\mathbf{\Sigma}$. Reflecting this, $h(\mathbf{t}) = \mathbf{\Sigma}\mathbf{V}^T \mathbf{t}$ is the functional confounder for an individual \mathbf{t} .

B.2 Expanded results

In [table 2](#), we list the 13 SNPs recovered by LODE, that have been previously reported as relevant to Celiac disease. In [fig. 7](#), we plot the true positive and false negative rate amongst SNPs deemed relevant by LODE. The ground truth here are the SNPs reported associated with celiac disease in prior literature.

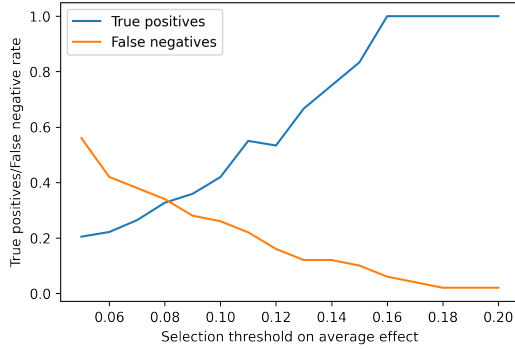


Figure 7: True positive vs. False negative rate as we vary the threshold on average effects, that determines which SNPs LODE deems relevant to the outcome.

SNP	EFFECT	LASSO COEF.
rs3748816	0.12	0.20
rs10903122	0.10	0.17
rs2816316	0.11	0.20
rs13151961	0.17	0.32
rs2237236	0.17	0.00
rs12928822	0.14	0.29
rs2187668	-0.70	-2.37
rs2327832	-0.12	-0.20
rs1738074	-0.16	-0.23
rs11221332	-0.15	-0.24
rs653178	-0.13	-0.21
rs4899260	-0.12	-0.19
rs17810546	-0.12	-0.20

Table 2: Full list of SNPs previously reported as relevant that were recovered by LODE, and their estimated effects and Lasso coefficients for SNPs. The effect threshold here is 0.1.