# Implicit Regularization in Overparameterized Bilevel Optimization

Paul Vicol[1,2], Jonathan Lorraine[1,2], David Duvenaud[1,2], Roger Grosse[1,2]
[1]University of Toronto, [2]Vector Institute

## Motivation & Summary

- Bilevel problems involve inner and outer parameters, each optimized for its own objective.

$$x^* \in \arg\min_x F(x, y^*)$$
$$y^* \in \mathcal{S}(x) = \arg\min_y f(x, y)$$

- **Examples**: hyperparameter optimization, dataset distillation, meta-learning, NAS, and GANs.
- Most prior work assumes that the inner & outer objectives have unique solutions, but often in practice, at least one of them is overparameterized → non-unique.
- We investigate the inductive biases of different gradient-based algorithms for jointly optimizing the inner and outer parameters.
- We distinguish between two different solution concepts—cold-start and warm-start equilibria
- The behavior depends on algorithmic choices such as the hypergradient approximation.

## Gradient-Based Bilevel Optimization

- Gradient-based bilevel opt requires the gradient of the outer objective with respect to the outer parameters, called the *hypergradient*. For a given solution $y^* \in \mathcal{S}(x)$, which is called a *best-response* to x:

$$\frac{dF(x, y^*)}{dx} = \frac{\partial F}{\partial x} + \frac{\partial F}{\partial y^*} \frac{\partial y^*}{\partial x}$$
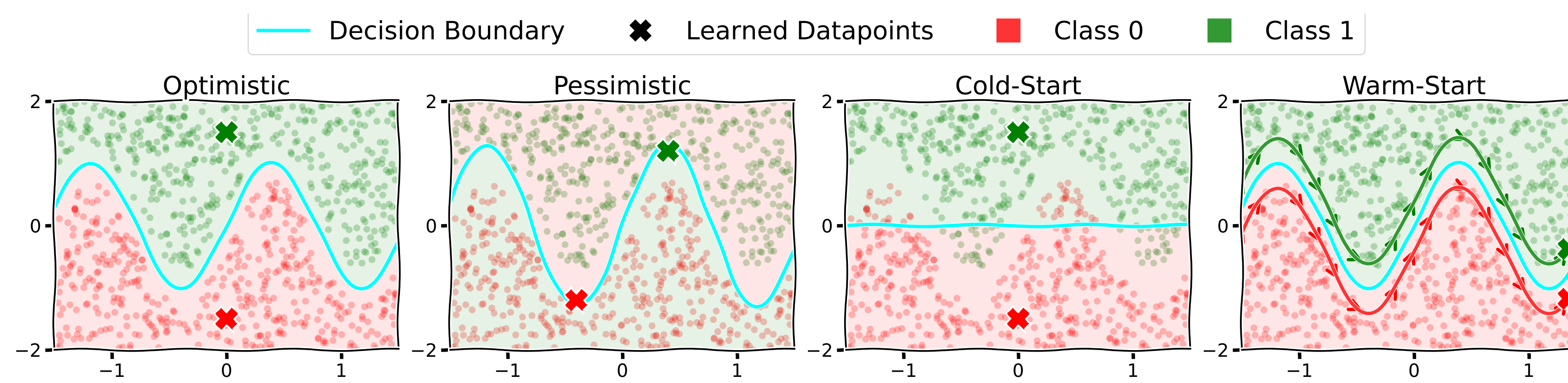
- Common ways to compute the response Jacobian are:
- Differentiation through unrolling: $\frac{dy^*}{dx} \approx \frac{d\Phi_k(y_0, x)}{dx}$
- Implicit differentiation: $\frac{dy^*}{dx} = -\left(\frac{\partial^2 f}{\partial y \partial y^\top}\right)^{-1} \frac{\partial^2 f}{\partial y \partial x}$
- Common approximations to the inverse Hessian include: 1) truncated CG, and 2) the truncated Neumann series:

$$\left(\frac{\partial^2 f}{\partial y \partial y^\top}\right)^{-1} \approx \sum_{j=0}^{K} \left(I - \frac{\partial^2 f}{\partial y \partial y^\top}\right)^j$$

## Warm-Start vs Cold-Start

- **Cold-start:** re-initialize the inner parameters and run the inner optimization to convergence each time we compute the gradient for the outer parameters
- **Warm-start:** jointly optimize the inner and outer parameters in an online fashion, e.g., alternating gradient steps with their respective objectives

## Warm-Start vs Cold-Start (Contd.)



Legend: Decision Boundary, Learned Datapoints (✗), Class 0 (red), Class 1 (green)
Panels: Optimistic, Pessimistic, Cold-Start, Warm-Start

### Solutions for Overparameterized Inner Problems

- The *optimistic* solution chooses the inner parameters that achieve the *best outer-objective value*, $\arg\min_{y \in \mathcal{S}(x)} F(x, y)$.
- The *pessimistic* solution chooses $y \in \mathcal{S}(x)$ that achieves the *worst outer-objective value*, $\arg\max_{y \in \mathcal{S}(x)} F(x, y)$.
- In practice, due to the implicit bias of gradient descent, the $y \in \mathcal{S}(x)$ we end up at depends on the inner initialization $y_0$: with cold-start, we obtain $y$ that minimize the distance from $y_0$: $\arg\min_{y \in \mathcal{S}(x)} ||y - y_0||_2^2$.
- With warm-start, the trajectory of outer parameters x during joint optimization (shown by the arrows) influences the inner parameters y.

### Proximal Inner Optimization

- We can formalize warm-started joint optimization by considering a proximally regularized inner objective: $y^* \in \arg\min_y \{f(x, y) + \frac{\epsilon}{2}||y - y_k||^2\}$

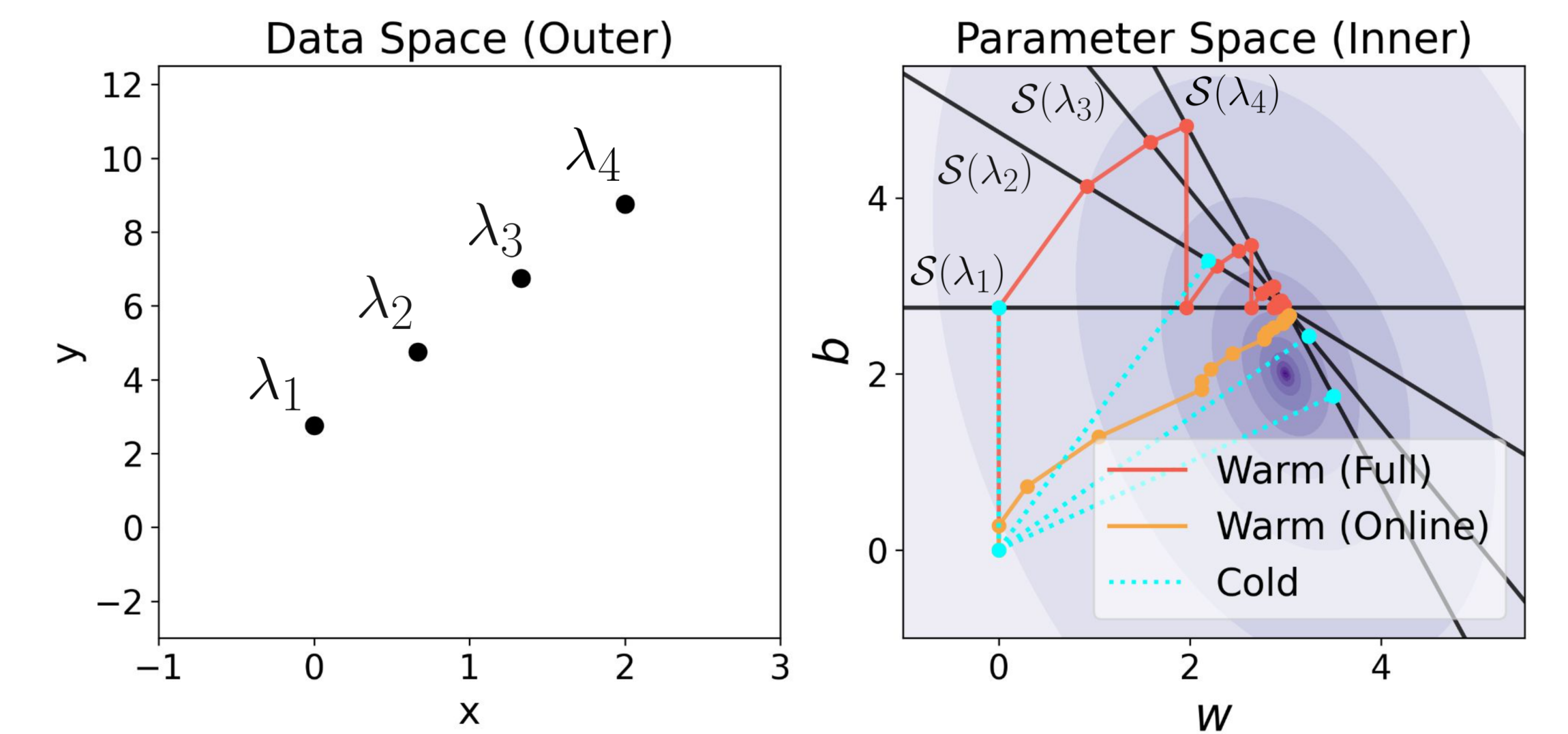| Cold-Start | Warm-Start |
|---|---|
| $x_{t+1} = x_t - \alpha \frac{\partial F}{\partial y^*} \frac{\partial y^*}{\partial x}$ | $x_{t+1} = x_t - \alpha \frac{\partial F}{\partial y_t^*} \frac{\partial y_t^*}{\partial x}$ |
| $y_{t+1}^* \in \arg\min_{y \in \mathcal{S}(x_{t+1})} ||y - y_0||^2$ | $y_{t+1}^* \in \arg\min_y \{f(x_{t+1}, y) + \frac{\epsilon}{2}||y - y_t||^2\}$ |

## Inner Overparameterization: Dataset Distillation



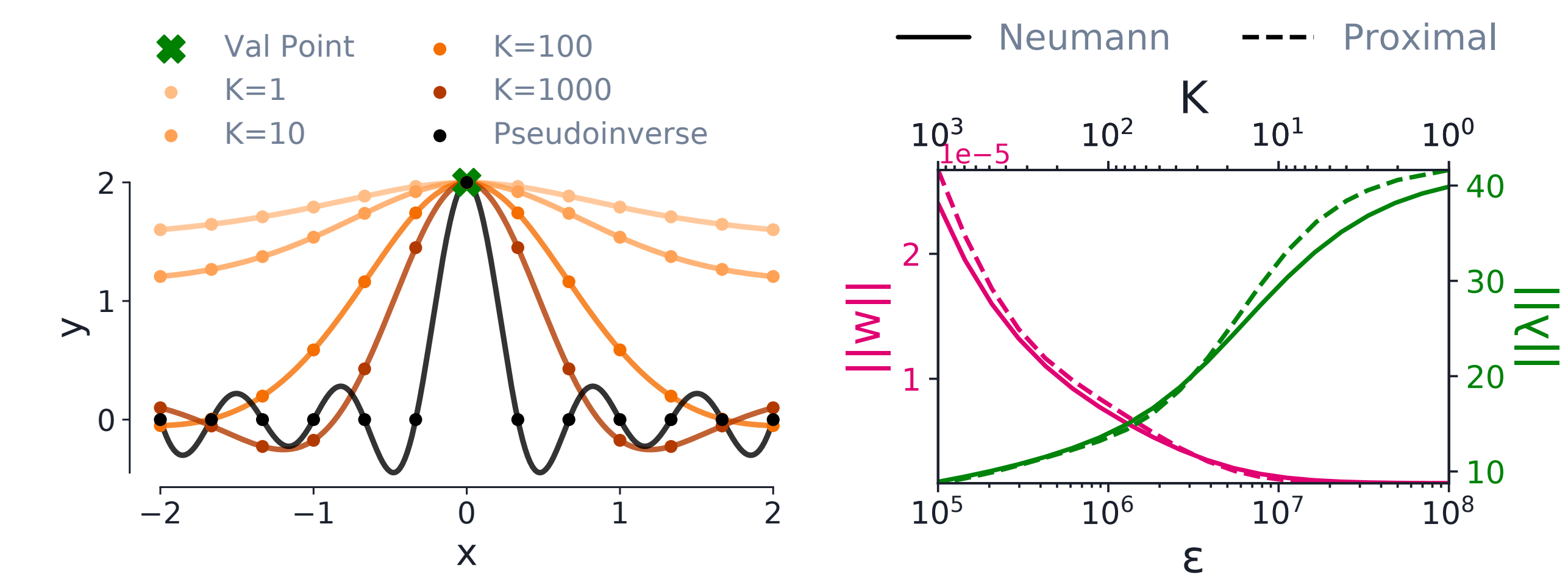Panels: Training on original data, Warm-start joint optim., Re-train on final points

- Dataset distillation for binary classification, with two learned datapoints (outer parameters) adapted jointly with the model weights (inner parameters).
- Because the outer obj is only used to update the outer params, one would think that all of the info about the outer obj is compressed into the outer params.
- Warm-starting yields a *trajectory* that traces out the boundary between classes.
- **Takeaway:** inner params can encode a surprising amount of information about the outer objective, even when the outer params are low-dimensional.

## Inner Overparameterization (Contd.)



Panels: Data Space (Outer), Parameter Space (Inner)
Legend: Warm (Full), Warm (Online), Cold

- Parameter-space view of warm-start with full inner optimization, warm-start with partial inner optimization (denoted "online"), and cold-start optimization.
- Cold-start projects from the origin onto the solution set for the current datapoint
- Warm-start projects from the current weights onto the solution set for the current datapoint
- By successive projection between solution sets, the inner parameters converge to the intersection of the solution sets, yielding inner params that perform well for multiple outer params simultaneously.

## Outer Overparameterization: Anti-Distillation



Legend: Val Point, K=1, K=10, K=100, K=1000, Pseudoinverse, Neumann, Proximal

- Fourier-basis 1D linear regression: we learn the y-coord of 13 synthetic datapoints such that a regressor trained on them will fit a single "val" datapoint, at the green X.
- **Left:** learned datapoints (outer params) from different hypergrad approximations: truncated Neumann/diff-through-unrolling with different # steps $K$
- **Right:** The norms of the inner and outer parameters, $||w - w_0||^2$ and $||\lambda - \lambda_0||^2$ as a function of $K$ (for Neumann/unrolling) or $\epsilon$ (for proximal).
- **Takeaway:** Empirically, the amount of inner optimization we perform affects the trade-off between the norms of the inner and outer params