

441 **Appendix**

442 **9 Preservation of Convexity and Submodularity**

443 **Proposition 1.** *If f is convex, then $g_P(\mathbf{y}, \theta) = f(P\mathbf{y}, \theta)$ is convex.*

444 *Proof.* The convexity can be simply verified by computing the second-order derivative:

$$\frac{d^2 g}{d\mathbf{y}^2} = \frac{d^2 f(P\mathbf{y}, \theta)}{d\mathbf{y}^2} = P^\top \frac{d^2 f}{d\mathbf{x}^2} P \succeq 0$$

445 where the last inequality comes from the convexity of f , i.e., $\frac{d^2 f}{d\mathbf{x}^2} \succeq 0$. □

446 **Proposition 2.** *If f is DR-submodular and $P \geq 0$, then $g_P(\mathbf{y}, \theta) = f(P\mathbf{y}, \theta)$ is DR-submodular.*

447 *Proof.* Assume f has the property of diminishing return submodularity (DR-submodular) [7]. Ac-
448 cording to definition of continuous DR-submodularity, we have:

$$\nabla_{\mathbf{x}_i, \mathbf{x}_j}^2 f(\mathbf{x}, \theta) \leq 0 \quad \forall i, j \in [n], \mathbf{y}$$

449 After applying the reparameterization, we can write:

$$g_P(\mathbf{y}, \theta) = f(\mathbf{x}, \theta)$$

450 and the second-order derivative:

$$\nabla_{\mathbf{y}}^2 g_P(\mathbf{y}, \theta) = P^\top \nabla_{\mathbf{x}}^2 f_P(\mathbf{x}, \theta) P \leq 0$$

451 Since all the entries of P are non-negative and all the entries of $\nabla_{\mathbf{x}}^2 f_P(\mathbf{x}, \theta)$ are non-positive by
452 DR-submodularity, the product $\nabla_{\mathbf{y}}^2 g_P(\mathbf{y}, \theta)$ also has all the entries being non-positive, which satisfies
453 the definition of DR-submodularity. □

454 **10 Quasiconvexity in Reparameterization Matrix**

455 **Proposition 3.** *$\text{OPT}(\theta, P) = \min_{\mathbf{y} \text{ feasible}} g_P(\mathbf{y}, \theta)$ is not globally quasiconvex in P .*

456 *Proof.* Without loss of generality, let us ignore the effect of θ and write $g_P(\mathbf{y}) = f(P\mathbf{x})$. In this
457 proof, we will construct a strongly convex function f where the induced optimal value function
458 $\text{OPT}(P) := \min_{\mathbf{y}} g_P(\mathbf{y})$ is not quasiconvex.

459 Consider $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3]^\top \in \mathbb{R}^3$. Define $f(\mathbf{x}) = \left\| \mathbf{x} - \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \right\|^2 \geq 0$ for all $\mathbf{x} \in \mathbb{R}^3$. Define $P =$

460 $\begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 2 \end{pmatrix}$ and $P' = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 2 & 0 \end{pmatrix}$. Apparently, $\mathbf{x}^* = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = P \begin{pmatrix} 1 \\ 0.5 \end{pmatrix}$ and $\mathbf{x}^* = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = P' \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$ are

461 both achievable. So the optimal values $\text{OPT}(P) = \text{OPT}(P') = 0$. But the combination $P'' = \frac{1}{2}P +$
462 $\frac{1}{2}P' = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \\ 1 & 1 \end{pmatrix}$ cannot, which results in an optimal value $\text{OPT}(P'') = \min_{\mathbf{y}} g_{P''}(\mathbf{y}) = > 0$

463 since $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \notin \text{span}(P'')$. This implies $\text{OPT}(\frac{1}{2}P + \frac{1}{2}P') = \text{OPT}(P'') > 0 = \frac{1}{2}\text{OPT}(P) + \frac{1}{2}\text{OPT}(P')$.

464 Thus $\text{OPT}(P)$ is not globally convex in the feasible domain. □

465 **Theorem 1.** *If $f(\cdot, \theta)$ is quasiconvex, then $\text{OPT}(\theta, P) = \min_{\mathbf{y} \text{ feasible}} g_P(\mathbf{y}, \theta)$ is quasiconvex in P_i
466 for any $1 \leq i \leq m$, where $P = [P_1, P_2, \dots, P_m] \geq 0$.*

467 *Proof.* Let us assume $P = [p_1, p_2, \dots, p_m]$ and $P' = [p'_1, p'_2, \dots, p'_m]$, where $p_i = p'_i \forall i \neq 1$ with
468 only the first column different. In the optimization problem parameterized by P , there is an optimal
469 solution $x = \sum_{i=1}^m p_i y_i$, $y_i \geq 0 \forall i$. Similarly, there is an optimal solution $x' = \sum_{i=1}^m p'_i y'_i$, $y'_i \geq 0 \forall i$
470 for the optimization problem parameterized by P' . We know that $f(x) = h(P)$, $f(x') = h(P')$.
471 Denote $P'' = cP + (1-c)P' = [p''_1, p''_2, \dots, p''_m]$ to be a convex combination of P and P' . Clearly,
472 $p''_1 = cp_1 + (1-c)p'_1$ and $p''_i = p_i = p'_i \forall i \neq 1$. Then we can construct a solution

$$\begin{aligned} x'' &= \frac{1}{\frac{c}{y_1} + \frac{1-c}{y'_1}} \left(\frac{c}{y_1} x + \frac{1-c}{y'_1} x' \right) \\ &= \frac{1}{\frac{c}{y_1} + \frac{1-c}{y'_1}} \left(\frac{c}{y_1} \sum_{i=1}^m p_i y_i + \frac{1-c}{y'_1} \sum_{i=1}^m p'_i y'_i \right) \\ &= \frac{1}{\frac{c}{y_1} + \frac{1-c}{y'_1}} (cp_1 + (1-c)p'_1) + \frac{1}{\frac{c}{y_1} + \frac{1-c}{y'_1}} \sum_{i=2}^m p_i \left(\frac{y_i}{y_1} + \frac{y'_i}{y'_1} \right) \\ &\in \text{Span}(P'') \end{aligned}$$

473 Thus, x'' is a feasible solution in the optimization problem parameterized by P'' . By the convexity of
474 f , we also know that

$$\begin{aligned} h(cP + (1-c)P') &= h(P'') \leq f(x'') \\ &= f\left(\frac{1}{\frac{c}{y_1} + \frac{1-c}{y'_1}} \left(\frac{c}{y_1} x + \frac{1-c}{y'_1} x' \right)\right) \\ &\leq \max(f(x), f(x')) \\ &= \max(h(P), h(P')) \end{aligned}$$

475 When one of y_1, y'_1 is 0, without loss of generality we assume $y_1 = 0$. Then we can construct
476 a solution $x'' = x$ which is still feasible in the optimization problem parameterized by $P'' =$
477 $cP + (1-c)P'$. Then we have the following:

$$h(P'') \leq f(x'') = f(x) = h(P) \leq \max(h(P), h(P'))$$

478 which concludes the proof. \square

479 11 Sample Complexity of Learning Predictive Model in Surrogate Problem

480 **Theorem 2.** Let \mathcal{H}_{lin} be the hypothesis class of all linear function mappings from $\xi \in \Xi \subset \mathbb{R}^p$ to
481 $\theta \in \Theta \in \mathbb{R}^n$, and let $P \in \mathbb{R}^{n \times m}$ be a linear reparameterization used to construct the surrogate. The
482 expected Rademacher complexity over t i.i.d. random samples drawn from \mathcal{D} can be bounded by:

$$\text{Rad}^t(\mathcal{H}_{lin}) \leq 2mC \sqrt{\frac{2p \log(2mt \|P^+\| \rho_2(S))}{t}} + O\left(\frac{1}{t}\right) \quad (4)$$

483 where C is the gap between the optimal solution quality and the worst solution quality, $\rho_2(S)$ is the
484 diameter of the set S , and P^+ is the pseudoinverse.

485 The proof of Theorem 2 relies on the results given by Balghithi et al. [11]. Balghithi et al. analyzed
486 the sample complexity of predict-then-optimize framework when the optimization problem is a
487 constrained linear optimization problem.

488 The sample complexity depends on the hypothesis class \mathcal{H} , mapping from the feature space Ξ to
489 the parameter space Θ . $\mathbf{x}_S^*(\theta) = \text{argmin}_{\mathbf{x} \in S} f(\mathbf{x}, \theta)$ characterizes the optimal solution with given
490 parameter $\theta \in \Theta$ and feasible region S . This can be obtained by solving any linear program solver
491 with given parameters θ . The optimization gap with given parameter P is defined as $\omega_S(\theta) :=$
492 $\max_{\mathbf{x} \in S} f(\mathbf{x}, \theta) - \min_{\mathbf{x} \in S} f(\mathbf{x}, \theta)$, and $\omega_S(\Theta) := \sup_{\theta \in \Theta} \omega_S(\theta)$ is defined as the upper bound on
493 optimization gap of all the possible parameter $\theta \in \Theta$. $\mathbf{x}^*(\mathcal{H}) := \{\xi \rightarrow \mathbf{x}^*(\Phi(\xi)) | \Phi \in \mathcal{H}\}$ is the
494 set of all function mappings from features ξ to the predictive parameters $\theta = \Phi(\xi)$ and then to the
495 optimal solution $\mathbf{x}^*(\theta)$.

496 **Definition 1** (Natarajan dimension). *Suppose that S is a polyhedron and \mathfrak{S} is the set of its extreme*
 497 *points. Let $\mathcal{F} \in \mathfrak{S}^\Xi$ be a hypothesis space of function mappings from Ξ to \mathfrak{S} , and let $A \in \Xi$ to be*
 498 *given. We say that \mathcal{F} shatters A if there exists $g_1, g_2 \in \mathcal{F}$ such that*

- 499 • $g_1(\xi) \neq g_2(\xi) \forall \xi \in A$.
- 500 • For all $B \subset A$, there exists $g \in \mathcal{F}$ such that (i) for all $\xi \in B, g(\xi) = g_1(\xi)$ and (ii) for all
- 501 $\xi \in A \setminus B, g(\xi) = g_2(\xi)$.

502 The Natarajan dimension of \mathcal{F} , denoted by $d_N(\mathcal{F})$, is the maximum cardinality of a set N -shattered
 503 by \mathcal{F} .

504 We first state their results below:

505 **Theorem 3** (Balghiti et al. [11] Theorem 2). *Suppose that S is a polyhedron and \mathfrak{S} is the set of its*
 506 *extreme points. Let \mathcal{H} be a family of functions mapping from features Ξ to parameters $\Theta \in \mathbb{R}^n$ with*
 507 *decision variable $\mathbf{x} \in \mathbb{R}^n$ and objective function $f(\mathbf{x}, \theta) = \theta^\top \mathbf{x}$. Then we have that*

$$\text{Rad}^t(\mathcal{H}) \leq \omega_S^*(\Theta) \sqrt{\frac{2d_N(\mathbf{x}^*(\mathcal{H})) \log(t|\mathfrak{S}|^2)}{t}}. \quad (5)$$

508 where Rad^t denotes the Radamacher complexity averaging over all the possible realization of t i.i.d.
 509 samples drawn from distribution \mathcal{D} .

510 The following corollary provided by Balghiti et al. [11] introduces a bound on Natarajan dimension
 511 of linear hypothesis class \mathcal{H} , mapping from $\Xi \in \mathbb{R}^p$ to $\Theta \in \mathbb{R}^n$:

512 **Corollary 1** (Balghiti et al. [11] Corollary 1). *Suppose that S is a polyhedron and \mathfrak{S} is the set of its*
 513 *extreme points. Let \mathcal{H}_{lin} be the hypothesis class of all linear functions, i.e., $\mathcal{H}_{\text{lin}} = \{\xi \rightarrow B\xi \mid B \in$
 514 $\mathbb{R}^{n \times p}\}$. Then we have*

$$d_N(\mathbf{x}^*(\mathcal{H}_{\text{lin}})) \leq np \quad (6)$$

515 Also $|\mathfrak{S}|$ can be estimated by constructing an ϵ -covering of the feasible region by open balls with
 516 radius ϵ . Let $\hat{\mathfrak{S}}_\epsilon$ be the centers of all these open balls. We can choose $\epsilon = \frac{1}{t}$ and the number of open
 517 balls required to cover S can be estimated by

$$|\hat{\mathfrak{S}}_\epsilon| \leq (2t\rho_2(S)\sqrt{n})^n \quad (7)$$

518 Combining Equation 5, 6, and 7, the Radamacher complexity can be bounded by:

Corollary 2 (Balghiti et al. [11] Corollary 2).

$$\text{Rad}^t(\mathcal{H}_{\text{lin}}) \leq 2n\omega_S(\Theta) \sqrt{\frac{2p \log(2nt\rho_2(S))}{t}} + O\left(\frac{1}{t}\right) \quad (8)$$

519 Now we are ready to prove Theorem 2:

520 *Proof of Theorem 2.* Now let us consider our case. We have a linear mapping from features $\xi \in$
 521 $X \subset \mathbb{R}^p$ to the parameters $\theta = B\xi \in \Theta \in \mathbb{R}^n$ with $B \in \mathbb{R}^{n \times p}$. The objective function is formed by

$$g_P(\mathbf{y}, \theta) = f(P\mathbf{y}, \theta) = \theta^\top P\mathbf{y} = (P^\top \theta)^\top \mathbf{y} = (P^\top B\xi)^\top \mathbf{y} \quad (9)$$

522 This is equivalent to have a linear mapping from $\xi \in \Xi \subset \mathbb{R}^p$ to $\theta' = P^\top B\xi$ where $P^\top B \in \mathbb{R}^{m \times p}$,
 523 and the objective function is just $g_P(\mathbf{y}, \theta') = \theta'^\top \mathbf{y}$. This yields a similar bound but with a smaller
 524 dimension $m \ll n$ as in Equation 10:

$$\text{Rad}^t(\mathcal{H}_{\text{lin}}) \leq 2m\omega_S(\Theta) \sqrt{\frac{2p \log(2mt\rho_2(S'))}{t}} + O\left(\frac{1}{t}\right) \quad (10)$$

525 where $\omega_S(\Theta)$ is unchanged because the optimality gap is not changed by the reparameterization. The
 526 only thing changed except for the substitution of m is that the feasible region S' is now defined in a

527 lower-dimensional space under reparameterization P . But since $\forall \mathbf{y} \in S'$, we have $P\mathbf{y} \in S$ too. So
 528 the diameter of the new feasible region can also be bounded by:

$$\begin{aligned}
 \rho(S') &= \max_{\mathbf{y}, \mathbf{y}' \in S'} \|\mathbf{y} - \mathbf{y}'\| \\
 &= \max_{\mathbf{y}, \mathbf{y}' \in S'} \|P^+ P(\mathbf{y} - \mathbf{y}')\| \\
 &= \max_{\mathbf{y}, \mathbf{y}' \in S'} \|P^+(P\mathbf{y} - P\mathbf{y}')\| \\
 &\leq \max_{\mathbf{x}, \mathbf{x}' \in S'} \|P^+(\mathbf{x} - \mathbf{x}')\| \\
 &\leq \|P^+\| \max_{\mathbf{x}, \mathbf{x}' \in S'} \|\mathbf{x} - \mathbf{x}'\| \\
 &= \|P^+\| \rho(S)
 \end{aligned}$$

529 where $P^+ \in \mathbb{R}^{m \times n}$ is the pseudoinverse of the reparameterization matrix P with $P^+P = I \in \mathbb{R}^{m \times m}$
 530 (assuming the matrix does not collapse). Substituting the term $\rho(S')$ in Equation 10, we can get the
 531 bound on the Radamacher complexity in Equation 4, which concludes the proof of Theorem 2. \square

532 12 Non-linear Reparameterization

533 The main reason that we use a linear reparameterization is to maintain the convexity of the inequality
 534 constraints and the linearity of the equality constraints. Instead, if we apply a convex reparameteriza-
 535 tion $\mathbf{x} = P(\mathbf{y})$, e.g., an input convex neural network [3], then the inequality constraints will remain
 536 convex but the equality constraints will no longer be affine anymore. So such convex reparameter-
 537 ization can be useful when there is no equality constraint. Lastly, we can still apply non-convex
 538 reparameterization but it can create non-convex inequality and equality constraints, which can be
 539 challenging to solve. All of these imply that the choice of reparameterization should depend on the
 540 type of optimization problem to make sure we do not lose the scalability while solving the surrogate
 541 problem.

542 13 Computing Infrastructure

543 All experiments were run on the computing cluster, where each node configured with 2 Intel Xeon
 544 Cascade Lake CPUs, 184 GB of RAM, and 70 GB of local scratch space. Within each experiment,
 545 we did not implement parallelization. So each experiment was purely run on a single CPU core. The
 546 main bottleneck of the computation is on solving the optimization problem, where we use Scipy [41]
 547 blackbox optimization solver. No GPU was used to train the neural network and throughout the
 548 experiments.