

# The Trade-offs of Incremental Linearization Algorithms for Nonsmooth Composite Problems

Krishna Pillutla, Vincent Roulet, Sham Kakade, Zaid Harchaoui



## Setting

Consider the finite sum composite problem

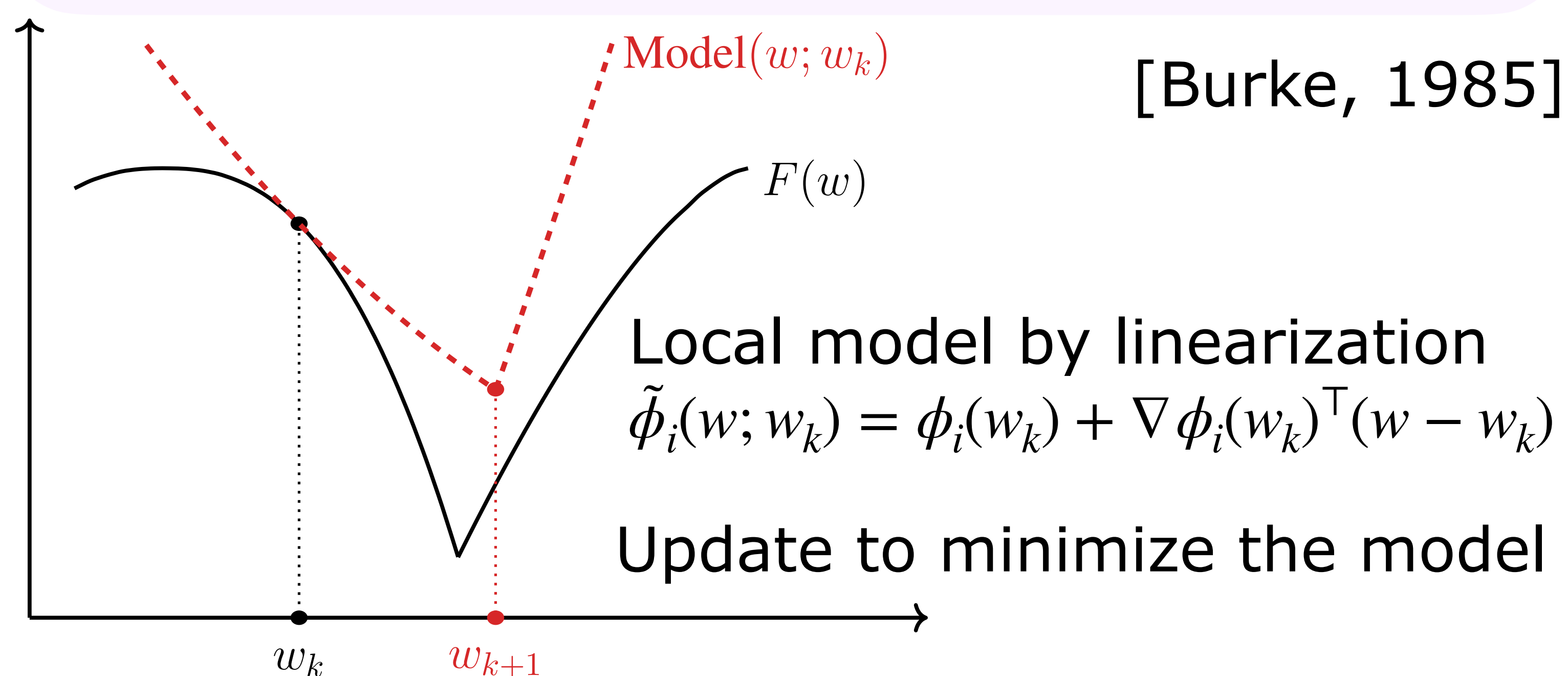
$$F(w) = \frac{1}{n} \sum_{i=1}^n f \circ \phi_i(w) \quad (1)$$

with  $f: \mathbb{R}^m \rightarrow \mathbb{R}$  Lipschitz and convex (loss),  
 $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}^m$  smooth and non-convex (predictor)

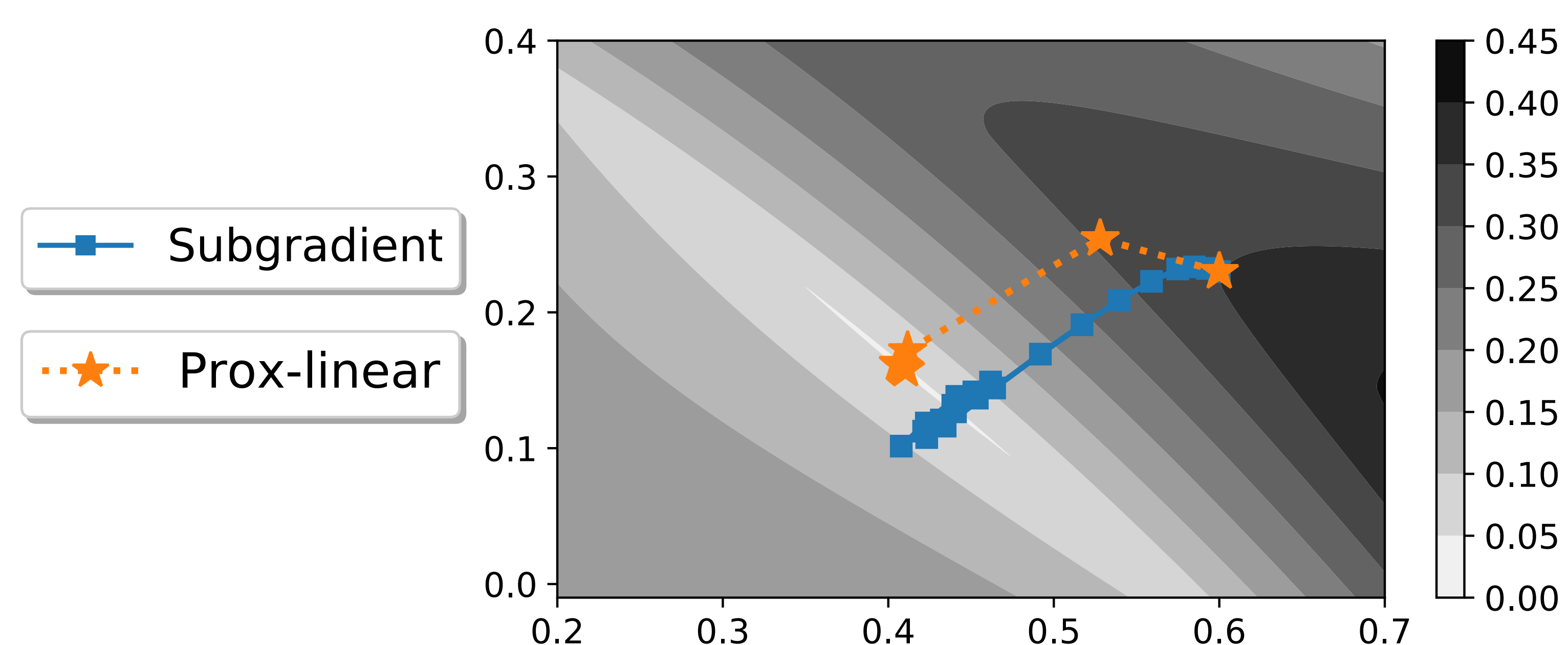
### Examples:

- Robust Regression:  $f = \|\cdot\|_2$
- Classification:  $f = \text{multi-class hinge loss}$

## Prox-linear/ Modified Gauss-Newton Method



$$w_{k+1} = \arg \min_w \frac{1}{n} \sum_{i=1}^n f(\tilde{\phi}_i(w; w_k)) + \frac{\kappa}{2} \|w - w_k\|_2^2$$



## Quadratic local convergence

**Proposition:** If  $f$  is  $\ell$ -Lipschitz and  $\mu$ -sharp,  $\phi = (\phi_1; \dots; \phi_n)$  is  $L$ -smooth and  $\sigma_{\min}(\nabla \phi(w)^\top) \geq \nu > 0$ , then  $F(w_k) \rightarrow F^*$  globally.

If  $F(w_k) - F^* \leq R$ , then for all  $t \geq k$ :

$$F(w_{t+1}) - F^* \leq \frac{1}{2R^2} (F(w_t) - F^*)^2$$

where  $R = \frac{\mu^2 \nu^2}{L \ell n^{3/2}}$

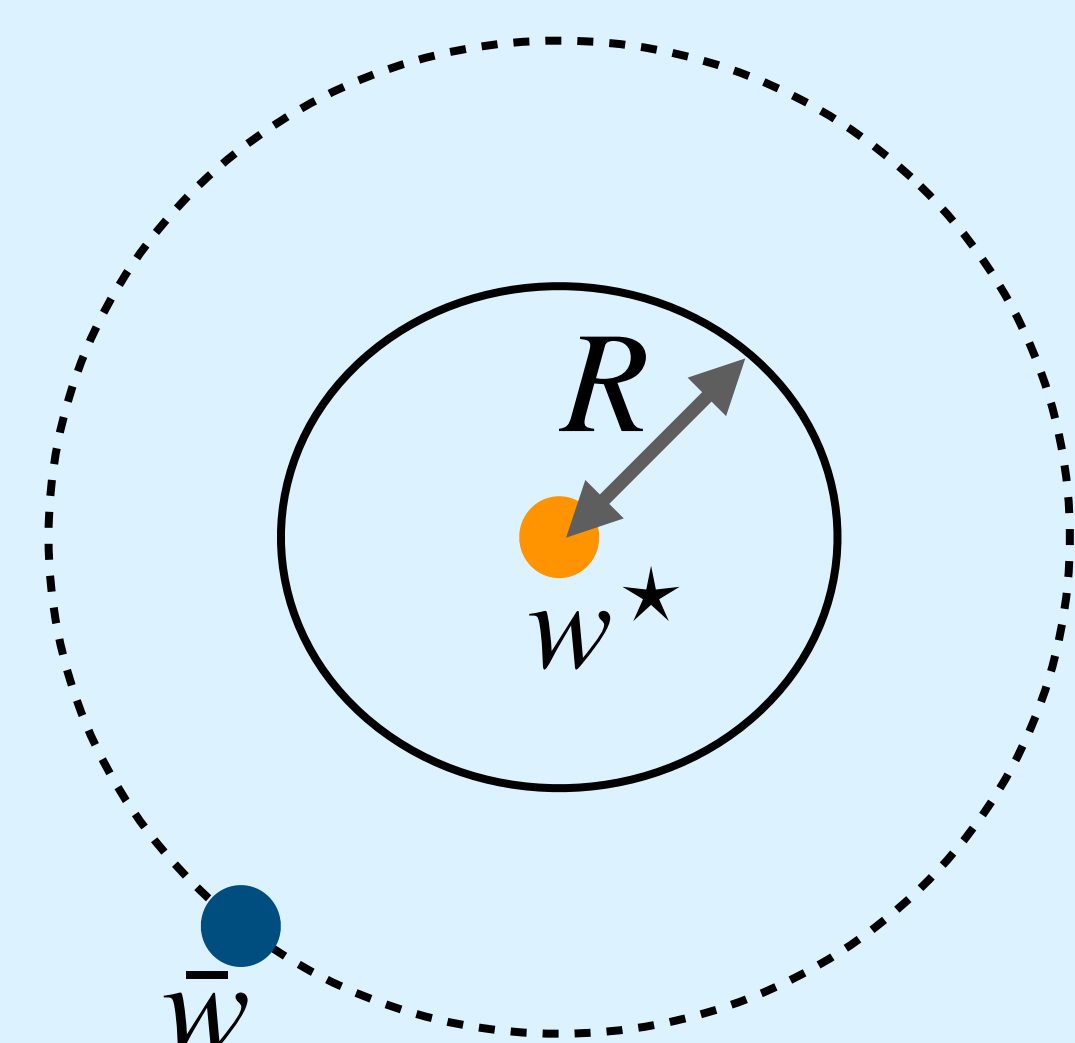
## Statistical Trade-offs

**Setting:**  $y_i = \psi(x_i; \bar{w}) + \xi_i$ , where  $\xi_i \sim \mathcal{N}(0, \sigma^2 I_m)$

Consider problem (1) with  $\phi_i(w) = \psi(x_i, w)$

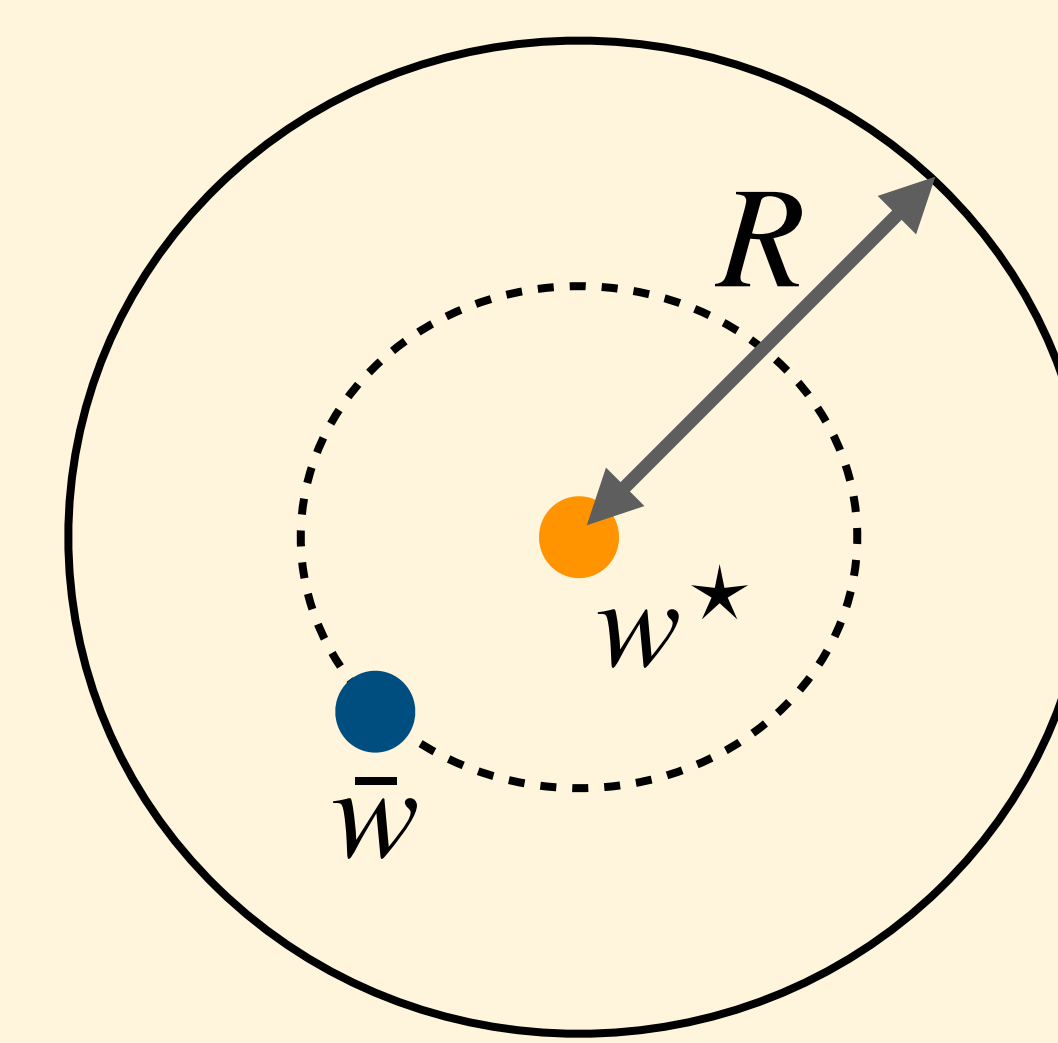
**Proposition:** Let  $R$  denote the radius of quadratic convergence and  $\exists w^*$  s.t.  $y_i = \psi(x_i, w^*)$  for all  $i$ . Let  $w_k$  be the first iterate enjoying quadratic convergence. Then w.h.p.,

(1) if  $\sigma > \frac{R}{m^{1/2} - m^{1/4}}$   
 then  $F(w_k) < F(\bar{w})$



$\sigma$  large  $\Rightarrow$  quadratic convergence is not active

(2) if  $\sigma < \frac{R}{m^{1/2} + m^{1/4}}$   
 then  $F(w_k) > F(\bar{w})$



$\sigma$  tiny  $\Rightarrow$  quadratic convergence is active

## Experiments

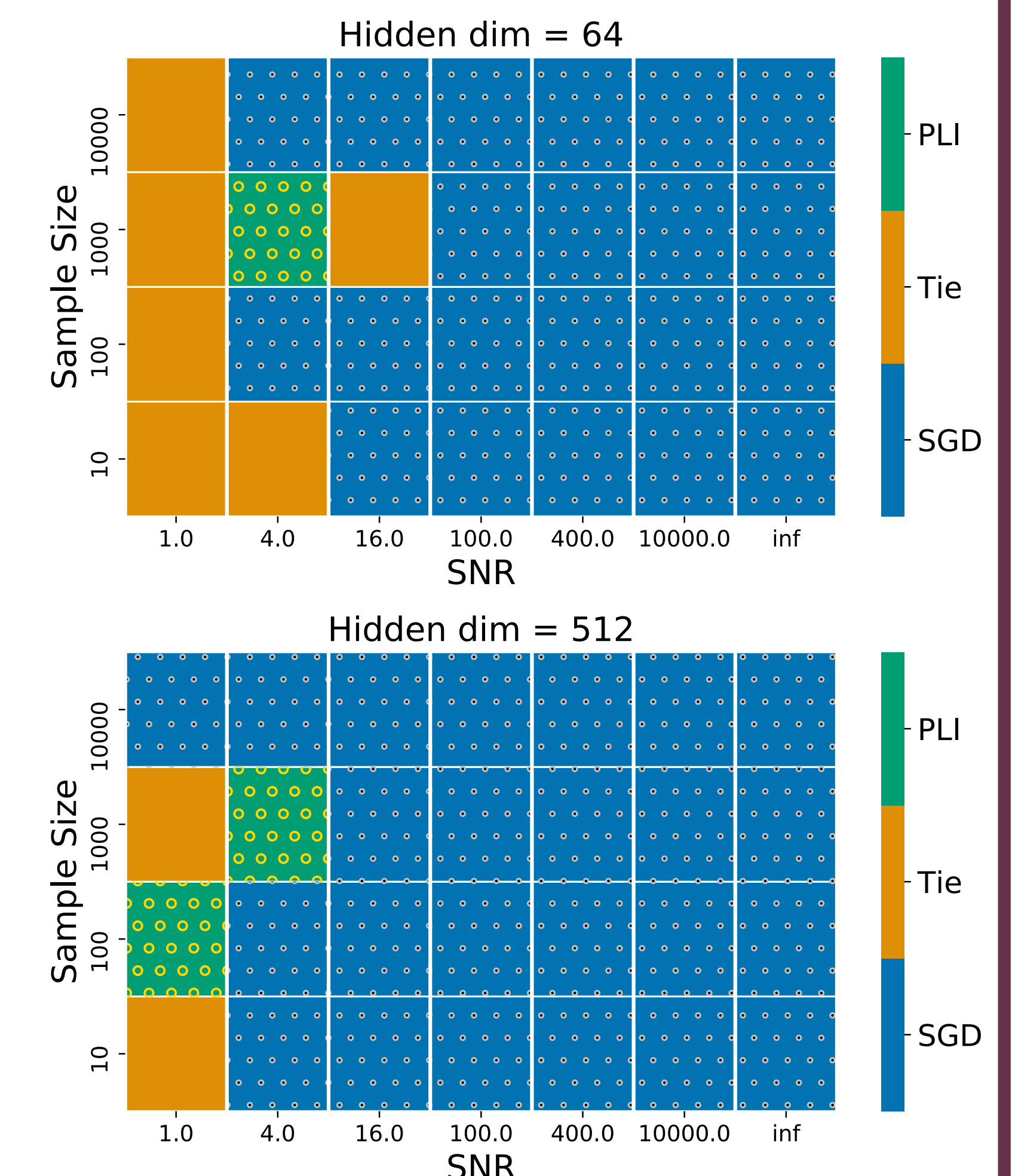
### Robust multi-output regression

**Input:**  $x \in \mathbb{R}^p$

**Output:**  $y \in \mathbb{R}^m$   
 (same as the statistical setting)

**Loss:**  $f = \|\cdot\|_2$

**Model:** 2-layer MLP



### Path planning (structured prediction)

**Input:** Image of a Warcraft map

**Output:** Least cost path from start to finish

**Loss:** structural hinge loss

**Model:** convolutional net

