



Competitive optimization

Two agents choose their decision variables to optimize their own objective. Conflicting objectives depend on both players' actions.

$$\min_x f(x, y), \quad \min_y g(x, y)$$

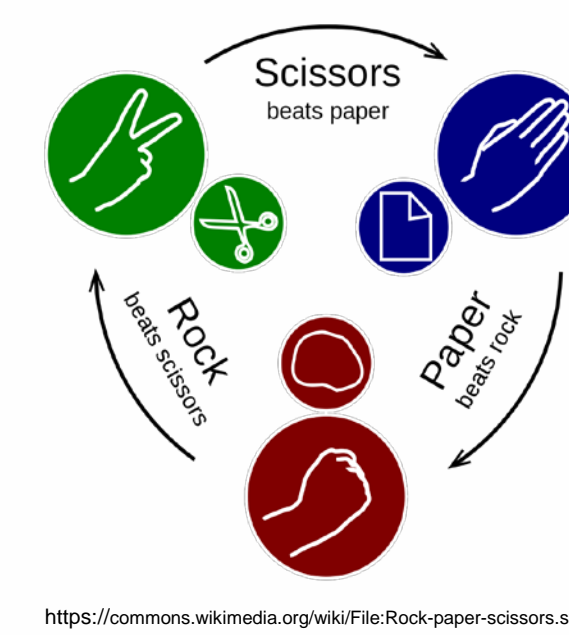
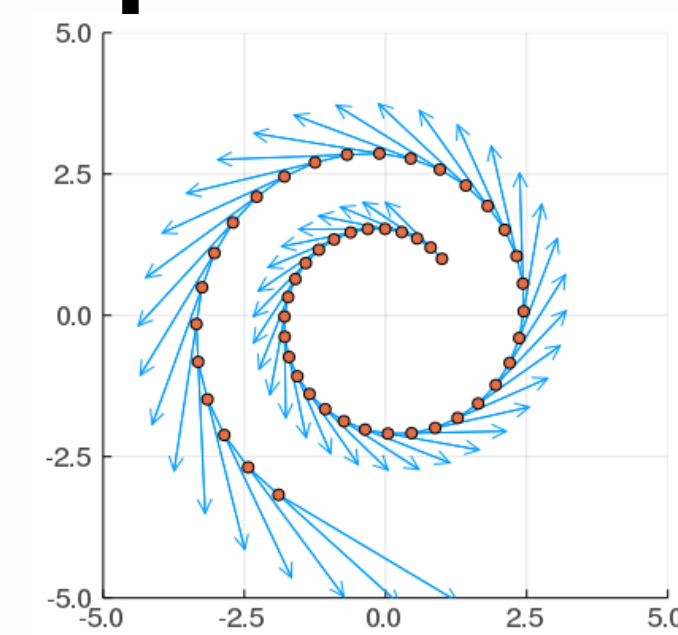
Arises in constrained optimization, robust statistics, and ML (GANs).

Rock! Paper! Scissor! Rock! Paper! ...

Naïve: Simultaneous Gradient Descent

$$\begin{aligned} x_{k+1} &= x_k - \eta \nabla_x f(x_k, y_k), \\ y_{k+1} &= y_k - \eta \nabla_y g(x_k, y_k) \end{aligned}$$

Divergent behavior even for simple bilinear game $f(x, y) = -g(x, y) = xy$!



What is gradient descent in two-player games?

Gradient descent minimizes quadratically regularized first order approximation:

$$x_{k+1} = x_k + \underset{x}{\operatorname{argmin}} f(x_k) + x^T \nabla_x f(x_k) + \frac{x^T x}{2\eta}$$

Thus, the generalization should be obtained from the **Nash equilibrium of a quadratically regularized first order approximation**

Linear or Bilinear approximation?

“Linear for one player \Rightarrow Linear for two players” loose the interactive aspect:

$$\begin{aligned} x_{k+1} &= x_k + \underset{x}{\operatorname{argmin}} x^T \nabla_x f + x^T D_{xy}^2 f y + y^T \nabla_y f + \frac{x^T x}{2\eta} \\ y_{k+1} &= y_k + \underset{y}{\operatorname{argmin}} x^T \nabla_x g + x^T D_{xy}^2 g y + y^T \nabla_y g + \frac{y^T y}{2\eta} \end{aligned}$$

All derivatives evaluated in (x_k, y_k)

“Linear for one player \Rightarrow **Bilinear for two players**” leads to interactive local game!

The local game has a unique Nash equilibrium

Theorem: The local game has a unique Nash equilibrium given by

$$\begin{aligned} x &= -\eta (\operatorname{Id} - \eta^2 D_{xy}^2 f D_{yx}^2 g)^{-1} (\nabla_x f - \eta D_{xy}^2 f \nabla_y g) \\ y &= -\eta (\operatorname{Id} - \eta^2 D_{yx}^2 g D_{xy}^2 f)^{-1} (\nabla_y g - \eta D_{yx}^2 g \nabla_x f) \end{aligned}$$

Novel algorithm uses local Nash as update rule

Algorithm: [Competitive Gradient Descent (CGD)]:

At each step, compute (x_{k+1}, y_{k+1}) from (x_k, y_k) as

$$\begin{aligned} x_{k+1} &= x_k - \eta (\operatorname{Id} - \eta^2 D_{xy}^2 f D_{yx}^2 g)^{-1} (\nabla_x f - \eta D_{xy}^2 f \nabla_y g) \\ y_{k+1} &= y_k - \eta (\operatorname{Id} - \eta^2 D_{yx}^2 g D_{xy}^2 f)^{-1} (\nabla_y g - \eta D_{yx}^2 g \nabla_x f) \end{aligned}$$

What I think that they think that I think ...

Can write CGD update as

$$\begin{pmatrix} x_{k+1} - x_k \\ y_{k+1} - y_k \end{pmatrix} = -\eta \begin{pmatrix} \operatorname{Id} & \eta D_{xy}^2 f \\ \eta D_{yx}^2 g & \operatorname{Id} \end{pmatrix}^{-1} \begin{pmatrix} \nabla_x f \\ \nabla_y g \end{pmatrix}$$

Neumann series: $(\operatorname{Id} - A)^{-1} = \sum_{k=0}^{\infty} A^k$

- First Term: Simultaneous gradient descent; Optimal if other player stays still
- Second Term: Optimal if the other player thinks that the other player stays still
- Third Term: Optimal if the other player thinks that the other player thinks ...

Why bilinear is the right notion of first order

The CGD update rule is similar to a regularized and damped Newton's method

$$\begin{pmatrix} x_{k+1} - x_k \\ y_{k+1} - y_k \end{pmatrix} = -\eta \begin{pmatrix} \operatorname{Id} + \eta D_{xx}^2 f & \eta D_{xy}^2 f \\ \eta D_{yx}^2 g & \operatorname{Id} + D_{yy}^2 g \end{pmatrix}^{-1} \begin{pmatrix} \nabla_x f \\ \nabla_y g \end{pmatrix}.$$

Why drop the diagonal block of the Hessian? We argue that **CGD is not a second order method, but the right way to do first order for competitive problems:**

Reason 1: Bilinear fully uses first order regularity

Most competitive optimization problems have objectives of the form

$$f(x, y) = \Phi(X(x), Y(y)),$$

for highly regular Φ possibly less regular $x \mapsto X(x)$, $y \mapsto Y(y)$ In this setting, mixed second derivatives are well behaved, as soon as gradients are well behaved.

Reason 2: It plays well with quadratic regularization

For non-convex f or g , Newton step can be local *worst* strategy of a player.

This leads to spurious attractors of the dynamics. Bilinear approximation is highest order that always leads to local best strategies when using quadratic regularization.

Reason 3: It leads to the right invariance:

First or second order approximation are invariant under linear transformations:

$$(A^{-1}x)^T \nabla_x f(A \cdot) = x^T \nabla_x f, \quad (A^{-1}x)^T D_{xx}^2 f(A \cdot) (A^{-1}x) = x^T D_{xx}^2 f x.$$

For the bilinear approximation only satisfies

$$\left(A^{-1} \begin{pmatrix} x \\ y \end{pmatrix} \right)^T \begin{pmatrix} 0 & D_{xy}^2 f(A \cdot) \\ D_{yx}^2 f(A \cdot) & 0 \end{pmatrix} \begin{pmatrix} A^{-1} x \\ A^{-1} y \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix}^T \begin{pmatrix} 0 & D_{xy}^2 f \\ D_{yx}^2 f & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

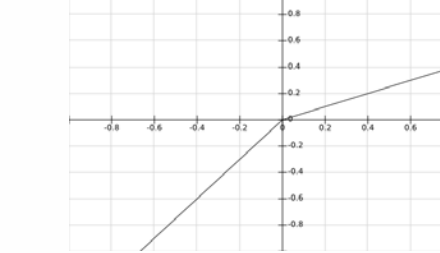
for $A = \begin{pmatrix} A_{xx} & 0 \\ 0 & A_{yy} \end{pmatrix}$ block diagonal. If A not block diagonal, it transfers decision variables between agents. This changes the game, don't want to be invariant to it!

Convergence results robust to strong interactions:

Consider zero sum game ($f = -g$) and define

$$\tilde{D} := (\operatorname{Id} + \eta^2 D_{xy}^2 f D_{yx}^2 f)^{-1} \eta^2 D_{xy}^2 f D_{yx}^2 f \quad \tilde{D} := (\operatorname{Id} + \eta^2 D_{yx}^2 f D_{xy}^2 f)^{-1} \eta^2 D_{yx}^2 f D_{xy}^2 f.$$

Define $h_{\pm}(\lambda) := \frac{\min(3\lambda, \lambda)}{2}$



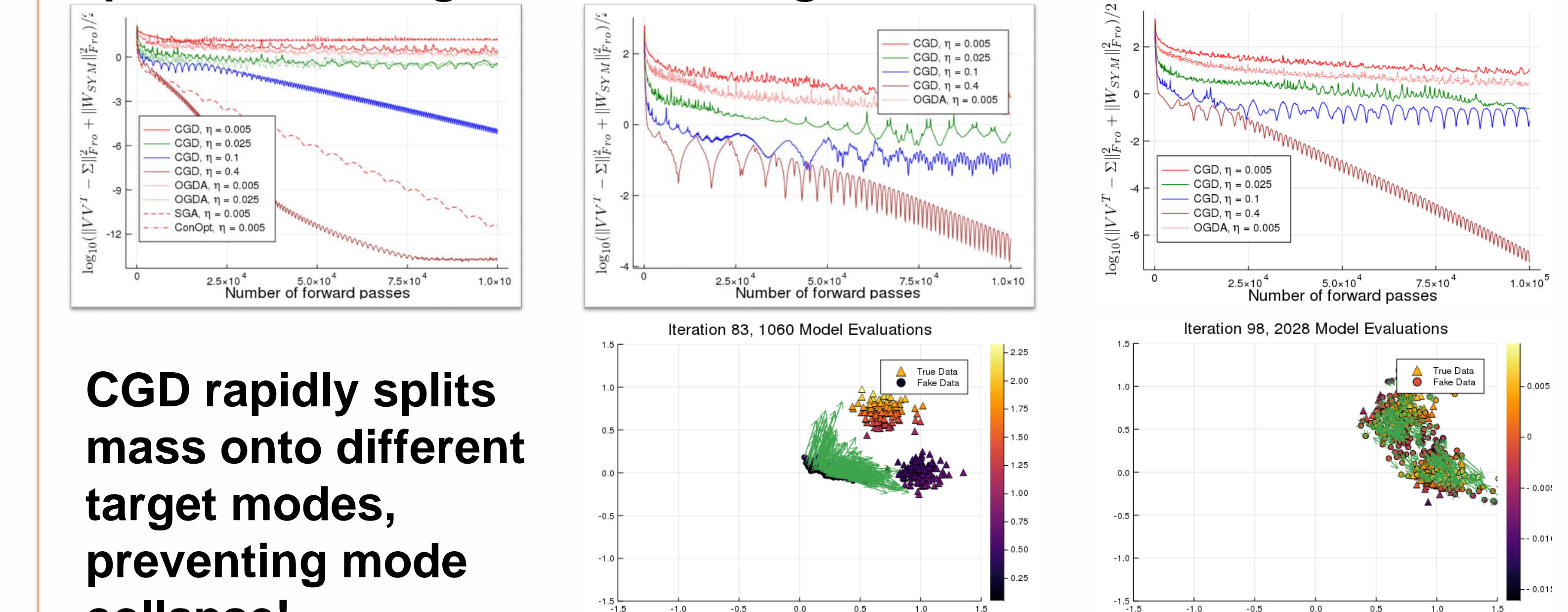
Theorem: If f is two times differentiable with L - Lipschitz continuous Hessian and $\eta \|D_{xx}^2 f\|, \eta \|D_{yy}^2 f\| \leq 1$:

$$\begin{aligned} & \| \nabla_x f(x_{k+1}, y_{k+1}) \|^2 + \| \nabla_y f(x_{k+1}, y_{k+1}) \|^2 - \| \nabla_x f \|^2 - \| \nabla_y f \|^2 \\ & \leq -\nabla_x f^T (\eta h_{\pm}(D_{xx}^2 f) + \tilde{D} - 32L\eta^2 \| \nabla_x f \|) \nabla_x f - \nabla_y f^T (\eta h_{\pm}(-D_{yy}^2 f) + \tilde{D} - 32L\eta^2 \| \nabla_y f \|) \nabla_y f \end{aligned}$$

Strong interaction between the players only improves convergence!

Faster convergence and splitting of modes

Improved convergence measuring number of model evaluations



CGD rapidly splits mass onto different target modes, preventing mode collapse!

Replacing gradient penalty by CGD improves WGAN-GP Inception Score on CIFAR10: Come to the poster session “Bridging Game Theory and Deep Learning”, Sat Dec14 9:30 am – 11 am and 4:30 pm in West Exhibition Hall A!