



Learning to Optimize

Ke Li Jitendra Malik

{ke.li, malik}@eecs.berkeley.edu

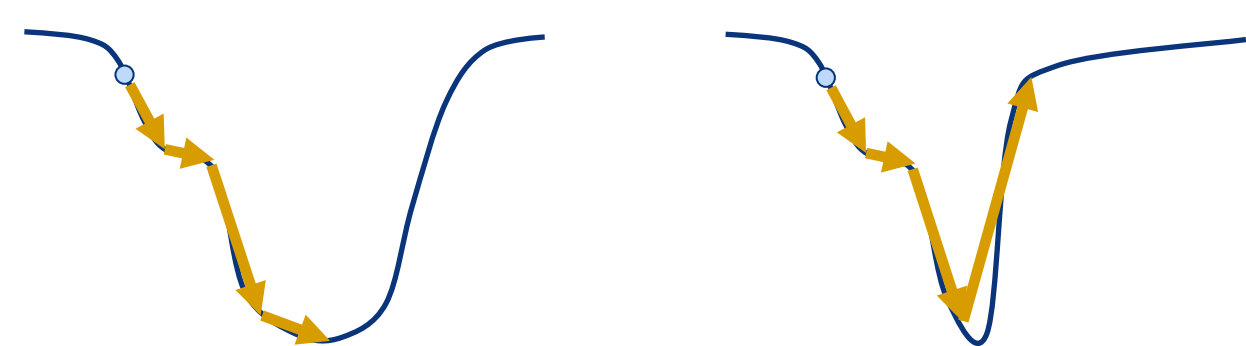
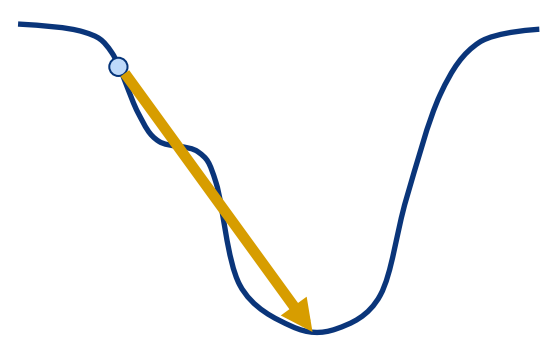
Introduction

- Optimization problems are ubiquitous in science and engineering.
- Devising a new optimization algorithm manually is challenging. Is there a better way?
- If the mantra of machine learning is to learn what is traditionally manually designed...

Why not learn the optimization algorithm itself?

Challenges

- This domain is prone to overfitting and underfitting.
- If we want to do well on a *single* objective function:
 - Consider an algorithm that memorizes the optimum.
 - This is the best optimizer since it gets to the optimum in one step.
- If we want to do well on *all* objective functions:
 - Given any optimizer, we can always construct an objective function that it performs poorly on.



- Goal: Do well on a class of objective functions with similar geometry, e.g.:
 - Logistic regression loss functions
 - Neural net classification loss functions

Setting

- Given: a set of training objective functions $f_1, \dots, f_n \sim \mathcal{F}$, a distribution \mathcal{D} for initializing the iterate and a meta-loss $\mathcal{L}(f, x^{(1)}, \dots, x^{(T)})$ that measures the quality of the iterates $x^{(1)}, \dots, x^{(T)}$.
- An optimization algorithm \mathcal{A} takes an objective function f and an initial iterate $x^{(0)}$ as input and produces a sequence of iterates $x^{(1)}, \dots, x^{(T)}$.
- Goal: learn \mathcal{A}^* such that $\mathbb{E}_{f \sim \mathcal{F}, x^{(0)} \sim \mathcal{D}} [\mathcal{L}(f, \mathcal{A}^*(f, x^{(0)}))]$ is minimized.
- We choose $\mathcal{L}(f, x^{(1)}, \dots, x^{(T)}) = \sum_{i=1}^T f(x^{(i)})$

Parameterizing Optimization Algorithms

Algorithm 1 General structure of optimization algorithms

```

Require: Objective function f
x^(0) ← random point in the domain of f
for i = 1, 2, ... do
  Δx ← φ(f, {x^(0), ..., x^(i-1)})
  if stopping condition is met then
    return x^(i-1)
  end if
  x^(i) ← x^(i-1) + Δx
end for

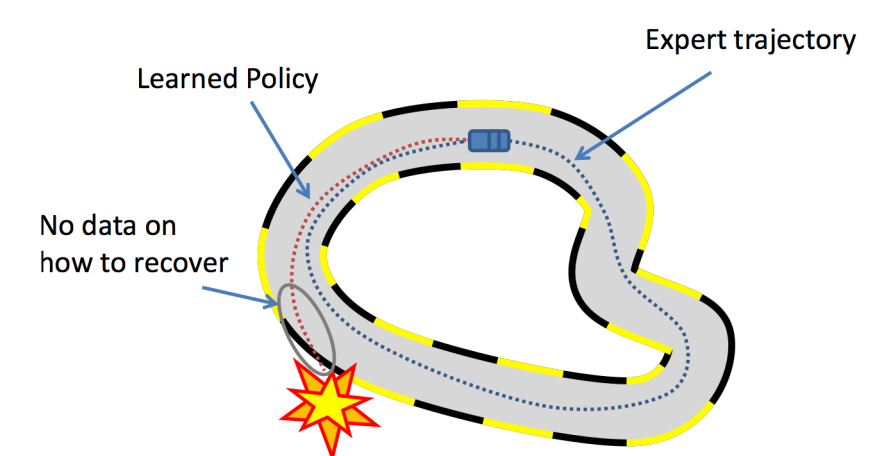
```

Gradient Descent $\phi(\cdot) = -\gamma \nabla f(x^{(i-1)})$
Momentum $\phi(\cdot) = -\gamma \left(\sum_{j=0}^{i-1} \alpha^{i-1-j} \nabla f(x^{(j)}) \right)$
Learned Algorithm $\phi(\cdot) = \text{Neural Net}$

- Input: Recent history of iterates, gradients and objective values
- Output: Step vector Δx
- Searching over the space of optimization algorithms reduces to learning the parameters of the neural net.

Properties of the Learning Problem

- The prediction of the neural net at any point in time affects the inputs that it sees in the future.
- This violates the i.i.d. assumption in supervised learning.
- Compounding errors: A policy trained using supervised learning does not know how to recover from previous mistakes.
- A supervised learner that makes a mistake with probability ϵ incurs a cumulative error of $O(\epsilon T^2)$, rather than $O(\epsilon T)$. (Ross and Bagnell, 2010)



Credit: John Schulman

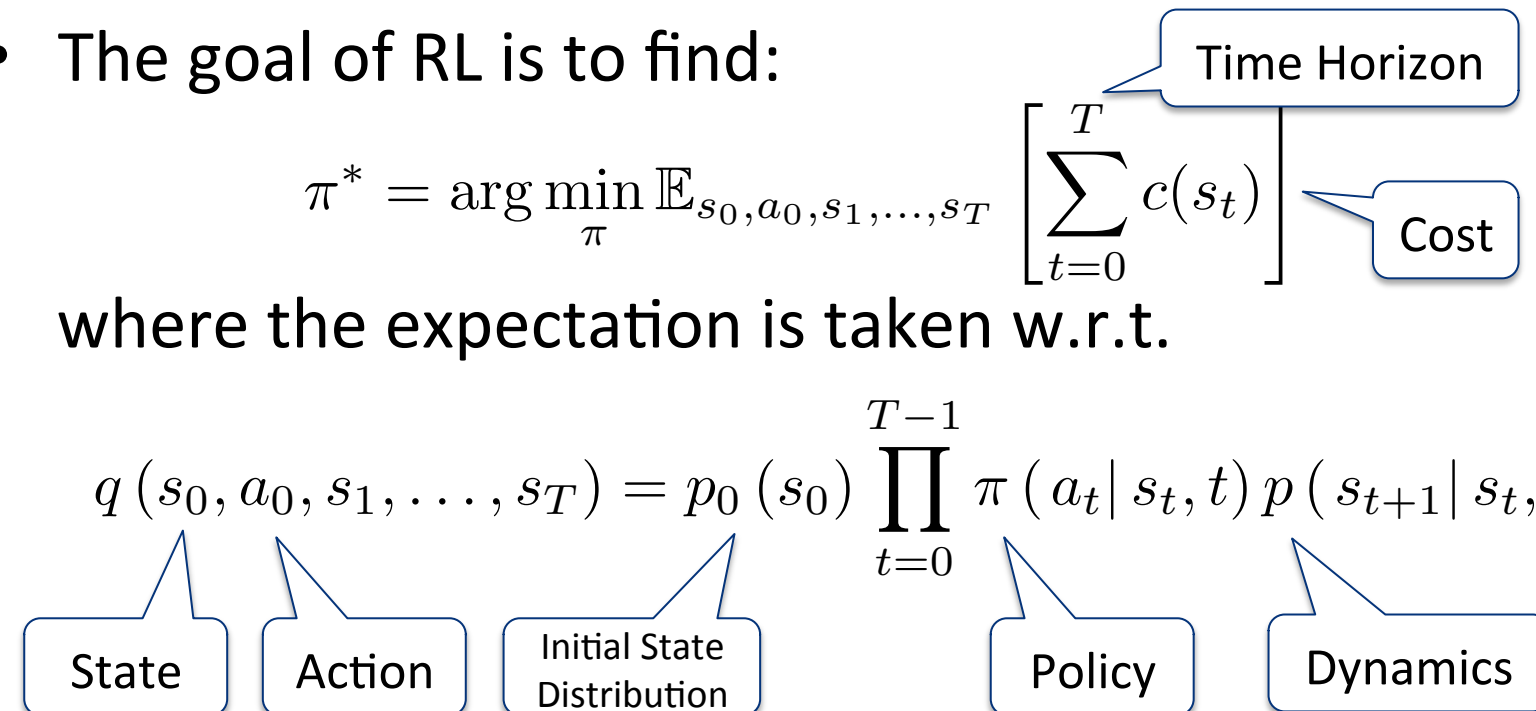
Reinforcement Learning

- The goal of RL is to find:

$$\pi^* = \arg \min_{\pi} \mathbb{E}_{s_0, a_0, s_1, \dots, s_T} \left[\sum_{t=0}^T c(s_t) \right]$$
 where the expectation is taken w.r.t.

$$q(s_0, a_0, s_1, \dots, s_T) = p_0(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t, t) p(s_{t+1} | s_t, a_t)$$
- The method we use is Guided Policy Search (Levine and Abbeel, 2014), which alternates between computing target trajectories and training the policy to replicate them. More precisely, it solves:

$$\min_{\theta, \eta} \mathbb{E}_{\psi} \left[\sum_{t=0}^T c(s_t) \right] \text{ s.t. } \psi(a_t | s_t, t; \eta) = \pi(a_t | s_t; \theta) \forall a_t, s_t, t$$



Formulation

Algorithm 1 General structure of optimization algorithms

```

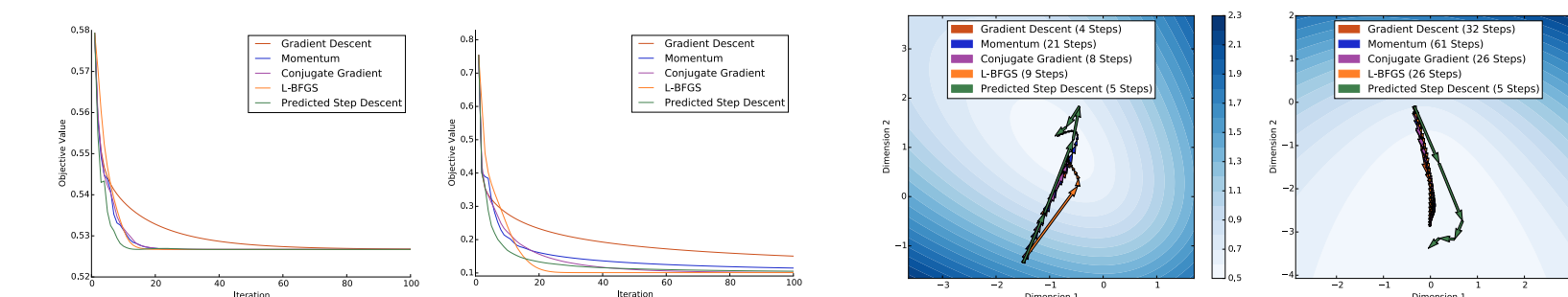
Require: Objective function f
x^(0) ← random point in the domain
for i = 1, 2, ... do
  Δx ← φ(f, {x^(0), ..., x^(i-1)})
  if stopping condition is met then
    return x^(i-1)
  end if
  x^(i) ← x^(i-1) + Δx
end for

```

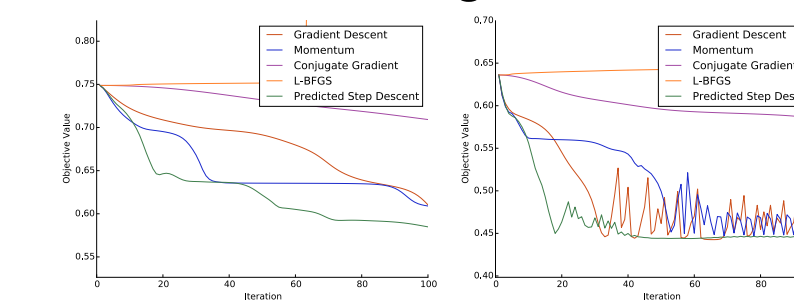
Experiments

- We trained optimizers for the following classes of low-dimensional optimization problems:
 - Logistic Regression (Convex)
 - Robust Linear Regression (Non-convex)
 - Small Neural Net Classifier (Non-convex)
- Trained on a set of random problems.
- Tested on a different set of random problems.

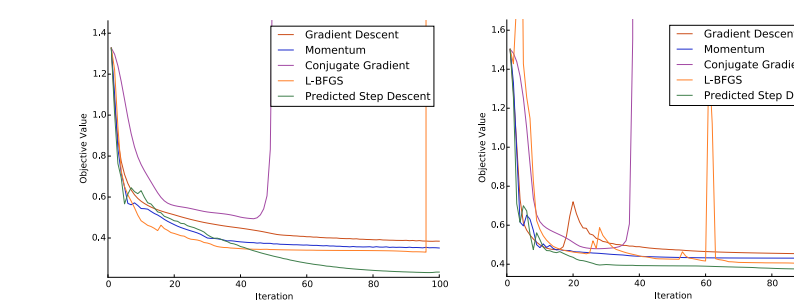
Logistic Regression:



Robust Linear Regression:



Small Neural Net:



arXiv:1606.01885 June 2016

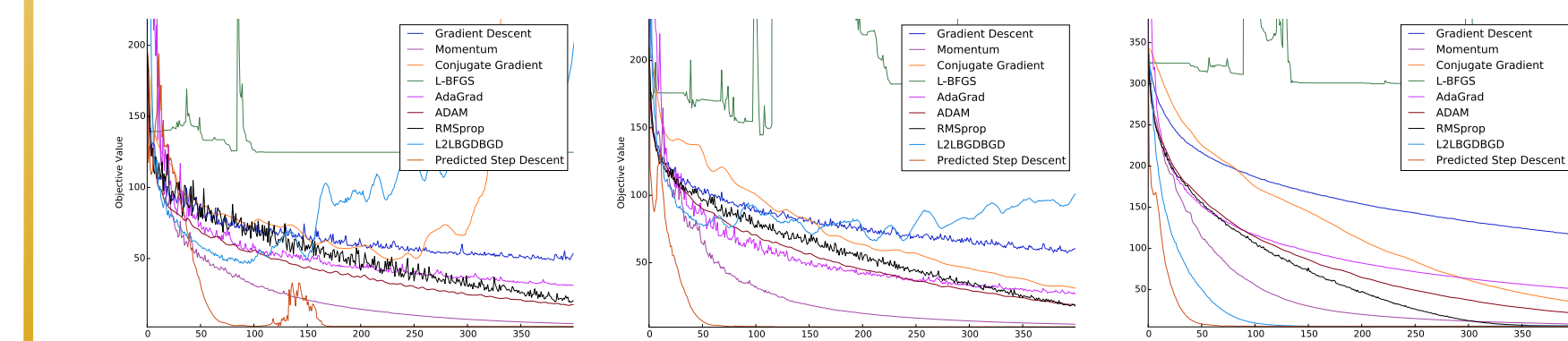
Future Work

Learning to Optimize Neural Nets

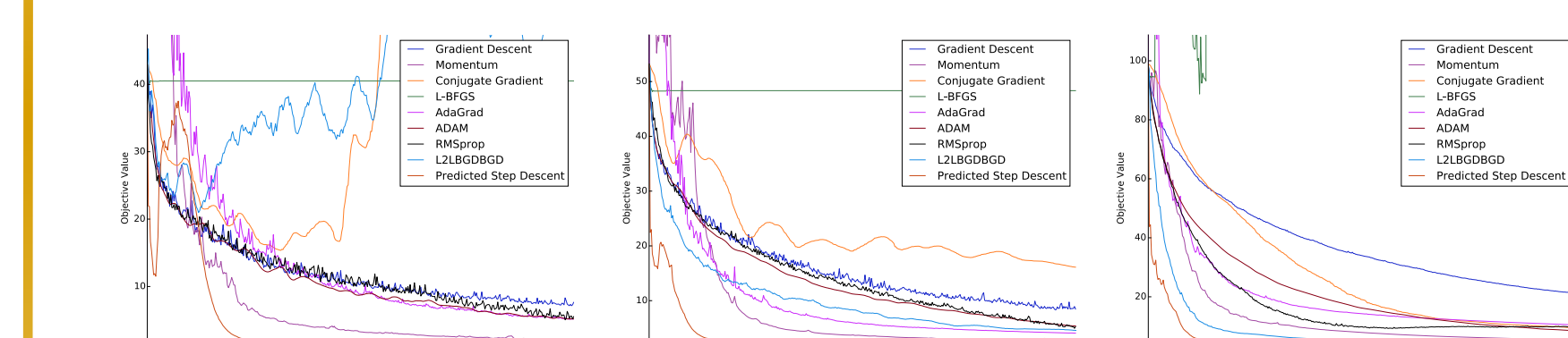
(<https://arxiv.org/abs/1703.00441>)

- Trained optimizer on the experience of training neural net on MNIST (a *single* objective function).
- Tested it on the problems of training a neural net on Toronto Faces, CIFAR-10 and CIFAR-100.

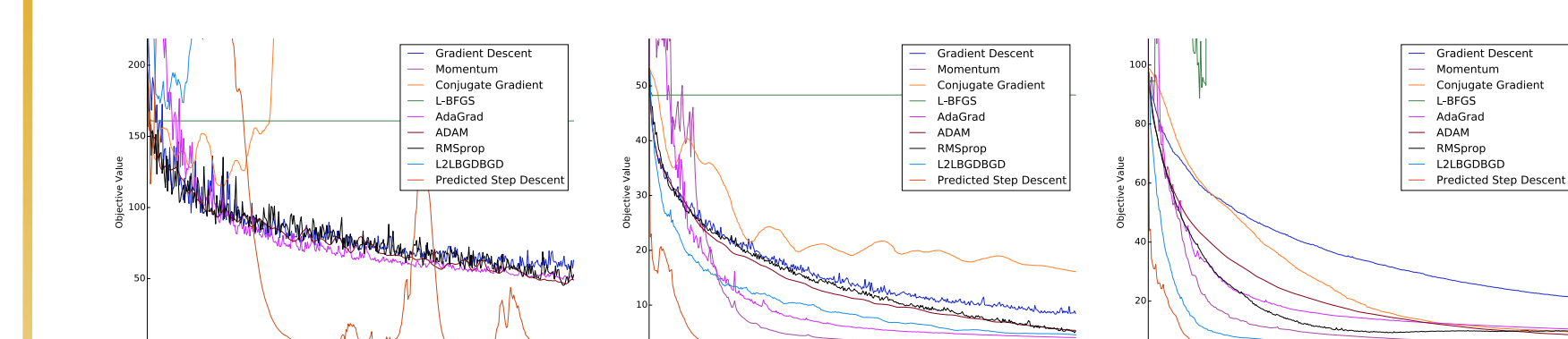
Wider Architecture:



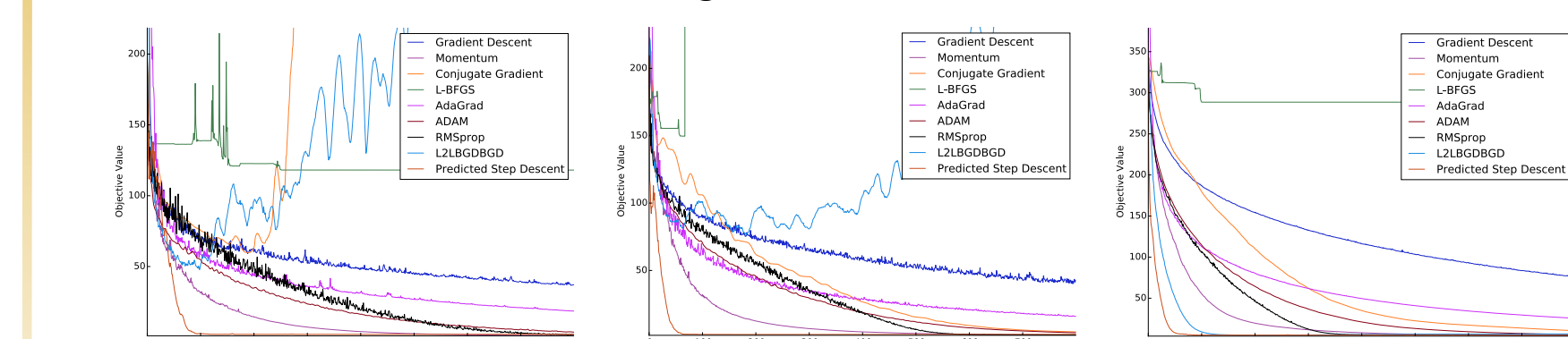
Noisier Gradients:



Wider Architecture and Noisier Gradients:



Wider Architecture and Longer Time Horizon:



arXiv:1703.00441 March 2017