

# Unconstrained Optimization

What is it?

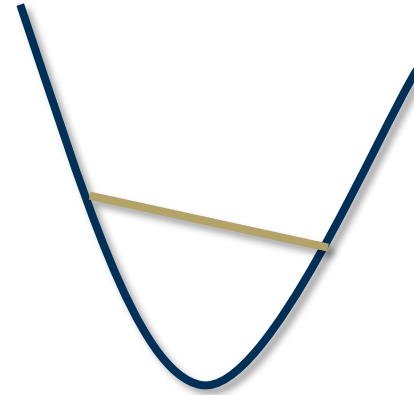
For our purposes today,

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x})$$

# The "Easy" Function Classes

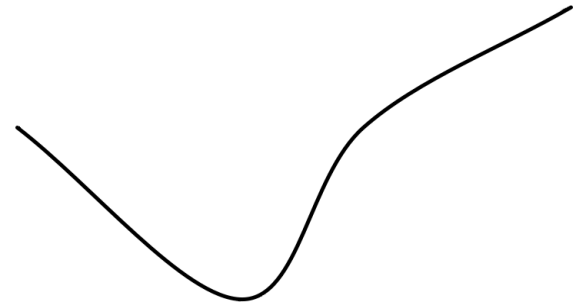
- Convex

$$f(\theta \mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y})$$



- PL (Polyak-Łojasiewicz)

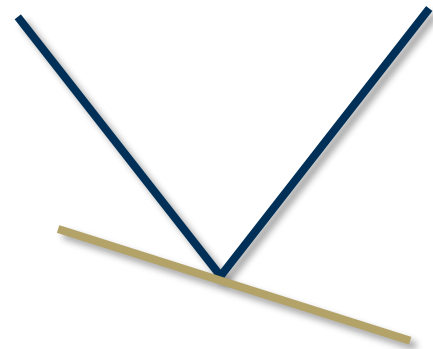
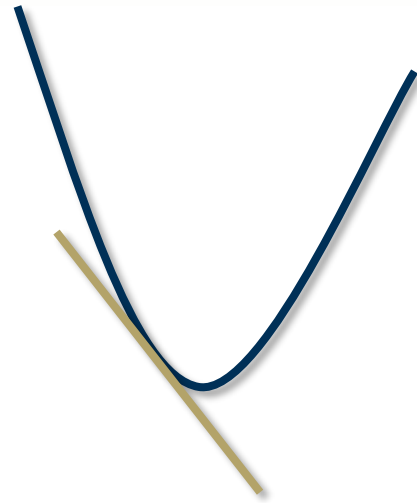
$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2m (f(\mathbf{x}) - f(\mathbf{x}^*))$$



(no spurious local minima)

# Non-Differentiable Convex Functions

- What is a gradient, really?
  - (for convex functions)
  - Linear lower bound
  
- Subgradient
  - Not necessarily unique at non-differentiable points
  - Subgradient analogs to gradient methods exist



# Regularization

Change the problem from

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x})$$

to

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \lambda r(\mathbf{x})$$

e.g.

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|^2$$

# Some Benefits of Regularization

- Better optimization landscape for chosen optimization algo
- Infinitely many solutions -> one unique solution
- Lower variance in solution / less overfitting
- Prefer certain solutions / incorporate domain knowledge
  - e.g. L1 regularizer promotes sparse solution

# Step-Size Line Search

- Sometimes gradient computations are expensive, but loss computations are cheap
- We need to make the most of every gradient computation we waited for

$$\underset{\alpha \geq 0}{\text{minimize}} \phi(\alpha) \quad \phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

- Sometimes this subproblem can be solved in closed form (e.g. convex quadratic problems)
- Otherwise, need a different approach...

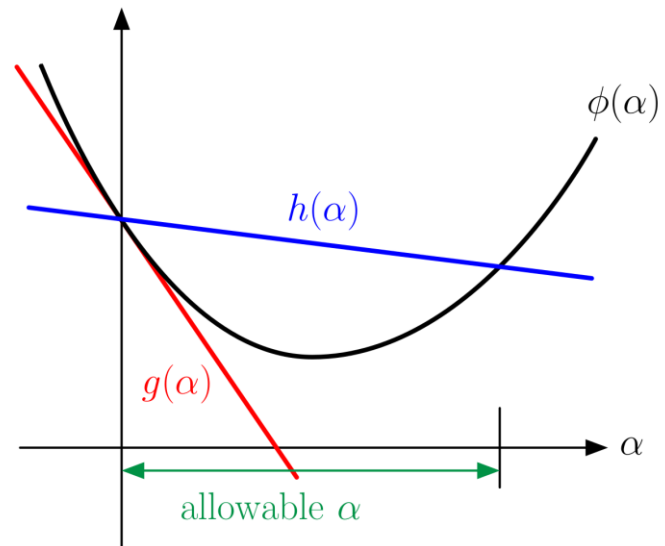
# Step-Size Line Search

- Backtracking: start with a big step, then shrink until *allowable*
  - *Allowable*: Armijo condition

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{d}_k) \geq c\alpha \langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle$$

$$\phi(\alpha) \leq h(\alpha) := f(\mathbf{x}_k) + c\alpha \langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle$$

$$\phi(\alpha) \geq g(\alpha) := f(\mathbf{x}_k) + \alpha \langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle$$



# Step-Size Line Search

## Backtracking line search

Input:  $\mathbf{x}_k$ ,  $\mathbf{d}_k$ ,  $\bar{\alpha} > 0$ ,  $c \in (0, 1)$ , and  $\rho \in (0, 1)$ .

Initialize:  $\alpha = \bar{\alpha}$

**while**  $\phi(\alpha) > h(\alpha)$  **do**

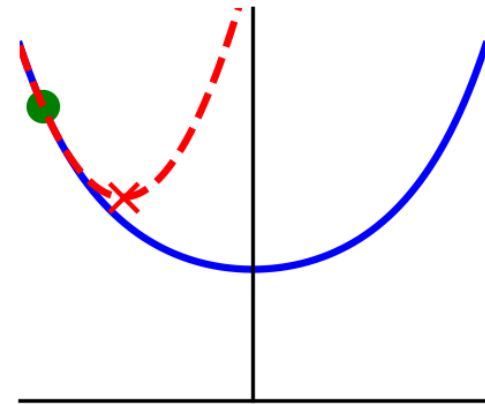
$\alpha = \rho\alpha$

**end while**

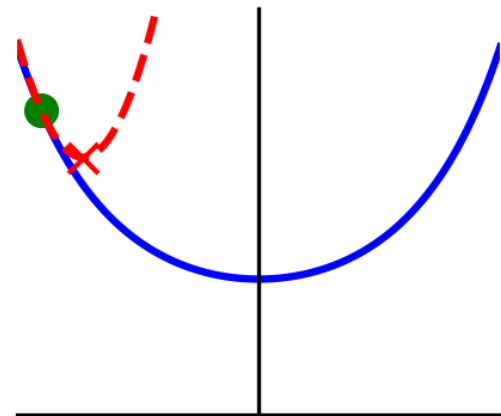


# Proximal Methods

Recall Newton's method uses local quadratic, which costs a Hessian computation



What if we just used *some* quadratic (that still agrees with the gradient)

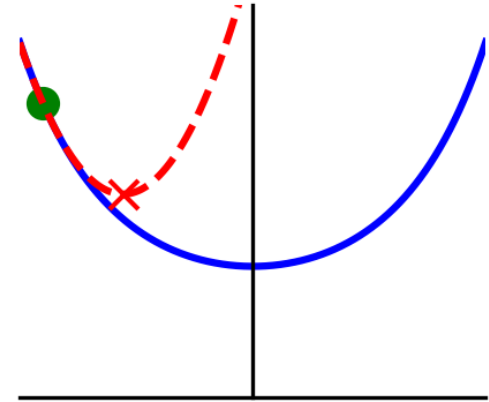


# Proximal Methods

$$\|\mathbf{x}\|_{\mathbf{A}}^2 := \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle$$

$$\underset{\mathbf{x}_{k+1}}{\text{minimize}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\nabla^2 f(\mathbf{x}_k)}^2$$

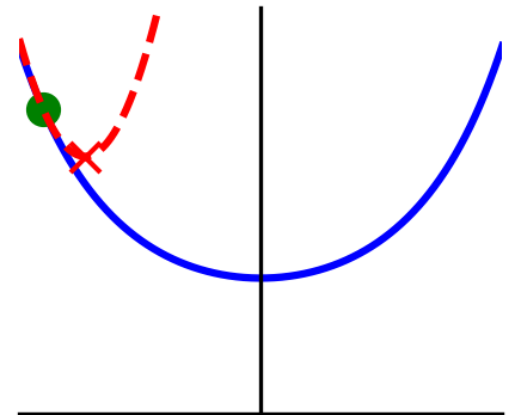
$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$



$$\underset{\mathbf{x}_{k+1}}{\text{minimize}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\frac{1}{\eta} \mathbf{I}}^2$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$$

(This is gradient descent!)



# Proximal Methods

The general form for a proximal step is

$$\underset{\mathbf{x}_{k+1}}{\text{minimize}} f(\mathbf{x}_k) + \frac{1}{2\eta} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$$

- To get GD, we substituted  $f(\mathbf{x})$  with its first-order approximation
- Sometimes it's better if we only use a first-order approximation for *part* of the objective function...

# Proximal Methods

The LASSO problem is

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1$$

What if we only used a first-order approximation for the L2 term, and kept the L1 term as-is? This is **ISTA** (iterative soft-thresholding algorithm), and it's way better than the standard subgradient method

$$f(\mathbf{x}_k) - f^*$$

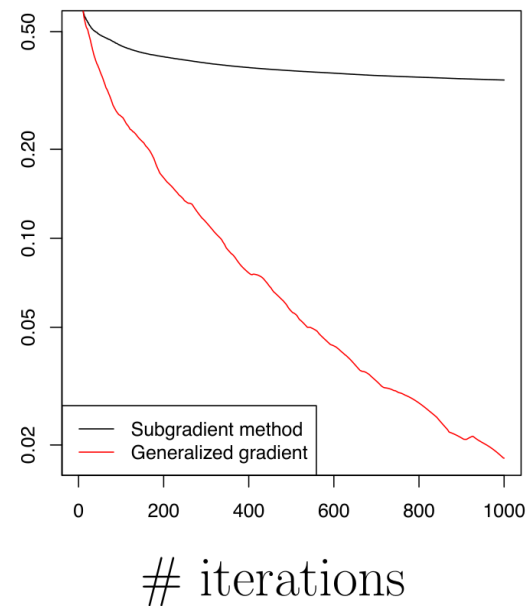


Chart taken from the lecture notes of Geoff Gordon and Ryan Tibshirani. ISTA is red.

# How Good is GD, Really?

- Hard to say in general, lots of things matter
- Strong Convexity
  - Lower bounded by (nonzero) quadratic

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}$$

- Smoothness
  - Upper bounded by quadratic

$$\nabla^2 f(\mathbf{x}) \preceq M\mathbf{I}$$

# How Good is GD, Really?

- When  $f$  is smooth, error decreases as  $O(1/k)$
- When  $f$  is smooth *and* strongly convex, error decreases as  $O(r^k)$  for some  $r < 1$  (called *linear convergence*)
- For reference, Newton's method has *quadratic convergence*, which is faster than linear convergence

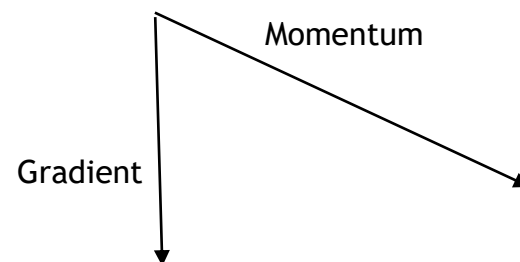
# GD with Momentum

- Some possible reasons to use momentum:
  - Better convergence rates (not all problems)
  - Better at avoiding spurious local minima
  - Prefers certain solutions to others (not always desirable)
  - Computing stochastic gradients (in batches), so want an integral-component to the steps

# GD with Momentum

- Polyak's heavy-ball method

```
for  $k = 0, 1, 2, \dots$  do
  if  $k = 0$  then
     $\mathbf{b}_{k+1} \leftarrow \nabla f(\mathbf{x}_k)$ 
  else
     $\mathbf{b}_{k+1} \leftarrow \mu \mathbf{b}_k + \nabla f(\mathbf{x}_k)$ 
  end if
   $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \gamma \mathbf{b}_{k+1}$ 
end for
```



- torch.optim.SGD
- $O(1/k)$  for smooth, strongly convex functions
- (Conjugate Gradients is a variant of heavy ball tuned for quadratics. It's the fast way to compute  $A^{-1}\mathbf{b}$  in high dimensions)



# GD with Momentum

- Nesterov's method

for  $k = 0, 1, 2, \dots$  do

  if  $k = 0$  then

$$\mathbf{b}_{k+1} \leftarrow \nabla f(\mathbf{x}_k)$$

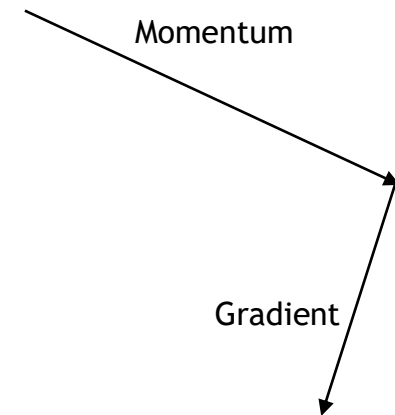
  else

$$\mathbf{b}_{k+1} \leftarrow \mu \mathbf{b}_k + \nabla f(\mathbf{x}_k)$$

  end if

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \gamma (\mu \mathbf{b}_{k+1} + \nabla f(\mathbf{x}_k))$$

end for



- torch.optim.SGD
- $O(1/k^2)$  for smooth, strongly convex functions

# GD with Momentum

- AdaGrad

$$\mathbf{s}_0 \leftarrow 0$$

**for**  $k = 0, 1, 2, \dots$  **do**

$$\mathbf{s}_{k+1} \leftarrow \mathbf{s}_k + (\nabla f(\mathbf{x}_k))^2 \quad (\text{element-wise})$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \gamma \frac{\nabla f(\mathbf{x}_k)}{\sqrt{\mathbf{s}_{k+1} + \epsilon}} \quad (\text{element-wise})$$

**end for**

- torch.optim.Adagrad
- More "equitable" trajectory
- Step size naturally shrinks over time

# GD with Momentum

- RMSprop

$$\mathbf{v}_0 \leftarrow 0$$

$$\mathbf{b}_0 \leftarrow 0$$

for  $k = 0, 1, 2, \dots$  do

$$\mathbf{v}_{k+1} \leftarrow \alpha \mathbf{v}_k + (1 - \alpha) (\nabla f(\mathbf{x}_k))^2 \quad (\text{element-wise})$$

$$\mathbf{b}_{k+1} \leftarrow \mu \mathbf{b}_k + \frac{\nabla f(\mathbf{x}_k)}{\sqrt{\mathbf{v}_{k+1} + \epsilon}} \quad (\text{element-wise})$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \gamma \mathbf{b}_{k+1}$$

end for

- torch.optim.RMSprop
- Gradient history is gradually "forgotten"
- Incorporates a momentum-like term

# GD with Momentum

- Adam

$$\mathbf{v}_0 \leftarrow 0$$

$$\mathbf{m}_0 \leftarrow 0$$

for  $k = 0, 1, 2, \dots$  do

$$\mathbf{v}_{k+1} \leftarrow \beta_2 \mathbf{v}_k + (1 - \beta_2) (\nabla f(\mathbf{x}_k))^2 \quad (\text{element-wise})$$

$$\mathbf{m}_{k+1} \leftarrow \beta_1 \mathbf{m}_k + (1 - \beta_1) \nabla f(\mathbf{x}_k) \quad (\text{element-wise})$$

$$\hat{\mathbf{v}}_{k+1} \leftarrow \frac{1}{1 - \beta_2^{k+1}} \mathbf{v}_{k+1}$$

$$\hat{\mathbf{m}}_{k+1} \leftarrow \frac{1}{1 - \beta_1^{k+1}} \mathbf{m}_{k+1}$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \gamma \frac{\hat{\mathbf{m}}_{k+1}}{\sqrt{\hat{\mathbf{v}}_{k+1} + \epsilon}} \quad (\text{element-wise})$$

end for

- torch.optim.Adam
- Uses a more classical momentum term
- A more dynamic approach increasing/decreasing step sizes

