

# OPTIMIZATION

# The right solution to a nearby problem? Or a nearby solution to the right problem?



- Chase the high-hanging fruit.
- Try to make stuff really work.
- Look for things that confuse/annoy you.



#### Alex Kipman

"The greatest danger for most of us is not that our aim is too high and we miss it, but that it is too low and we reach it" -- Michelangelo

Yesterday at 06:38 · Like · Comment

Write about it well.

### HOW TO GET IDEAS



- Computing PCA, LDA, ....
- Clustering, Gaussian mixture fitting, ....

Sometimes the optimization is easy...

...sometimes it's hard

Sometimes the hard problem is the one you must solve

# COMPUTER VISION AS OPTIMIZATION



#### Given function

# $f(x): \mathbb{R}^d \mapsto \mathbb{R}$ ,

# Devise strategies for finding x which minimizes f















>> print -dmeta







>> print –dpdf % then go to pdf and paste







>> set(gcf, 'paperUnits', 'centimeters', 'paperposition', [1 1 9 6.6]) >> print –dpdf % then go to pdf and paste

#### EXAMPLE





switch to matlab...





#### **ALTERNATION**

**Microsoft**<sup>®</sup>

- Ignore the hard ones
  - If you search, you'll find plenty of easy examples
  - But your customers won't...

• Or update more than one coordinate at once







- Alternation is slow because valleys may not be axis aligned
- So try gradient descent?

# **GRADIENT DESCENT**





- Alternation is slow because valleys may not be axis aligned
- So try gradient descent?

Steepest descent ( $x_0 = [0, 14]$ )

# GRADIENT DESCENT





- Alternation is slow because valleys may not be axis aligned
- So try gradient descent?

- Note that convergence proofs are available for both of the above
- But so what?

# GRADIENT DESCENT





#### AND ON A HARD PROBLEM

**Microsoft** 



- (Nonlinear) conjugate gradients
- Uses 1<sup>st</sup> derivatives only
- Avoids "undoing" previous work

SCG Scoled Conj Gran.

USE A BETTER ALGORITHM

**Microsoft**<sup>®</sup>



- (Nonlinear) conjugate gradients
- Uses 1<sup>st</sup> derivatives only
- And avoids "undoing" previous work
- 101 iterations on this problem

### USE A BETTER ALGORITHM





but we can do better...



- Starting with x how can I choose  $\delta$ so that  $f(x + \delta)$  is better than f(x)?
- So compute

$$\min_{\delta} f(x+\delta) \qquad g(\delta) = f(x+\delta)$$

 But hang on, that's the same problem we were trying to solve?

#### USE SECOND DERIVATIVES...



- Starting with x how can I choose  $\delta$ so that  $f(x + \delta)$  is better than f(x)?
- So compute  $\min_{\delta} f(x + \delta)$   $\approx \min_{\delta} f(x) + \delta^{\top} g(x) + \frac{1}{2} \delta^{\top} H(x) \delta$   $g(x) = \nabla f(x)$   $H(x) = \nabla \nabla^{\top} f(x)$

USE SECOND DERIVATIVES...





How does it look?

 $f(x) + \delta^{\mathsf{T}}g(x) + \frac{1}{2}\delta^{\mathsf{T}}H(x)\delta$  $g(x) = \nabla f(x) \in \mathbb{R}^{h}$   $H(x) = \nabla (\nabla^{T} f(x)) \in \mathbb{R}^{n \times n}$   $\int_{\mathcal{N}_{i}}^{\mathcal{N}_{i}} \int_{\mathcal{N}_{i}}^{\mathcal{N}_{i}} \int_{\mathcal{N}_{i$ drigx;

#### USE SECOND DERIVATIVES...

#### **Microsoft**



- Choose  $\delta$  so that  $f(x + \delta)$  is better than f(x)?
- Compute

 $\min_{\boldsymbol{\delta}} f + \boldsymbol{\delta}^{\mathsf{T}} g + \frac{1}{2} \boldsymbol{\delta}^{\mathsf{T}} H \boldsymbol{\delta}$ 

9+HS=0 HS=-9  $\delta = -H^{-1}g$  $= -H \int g$  || solve (H,g)

USE SECOND DERIVATIVES...



demo\_taylor\_2d(o, 'newton', 'rosenbrock')
demo\_taylor\_2d(1, 'newton', 'rosenbrock')
demo\_taylor\_2d(1, 'newton', `sqrt\_rosenbrock')

#### IS THAT A GOOD IDEA?



- Choose  $\delta$  so that  $f(x + \delta)$  is better than f(x)?
- Updates:

$$\boldsymbol{\delta}_{\text{Newton}} = -H^{-1}g$$

$$\boldsymbol{\delta}_{\text{GradientDescent}} = -\lambda g$$

#### USE SECOND DERIVATIVES...



• Updates:

$$\boldsymbol{\delta}_{\text{Newton}} = -H^{-1}g$$
$$\boldsymbol{\delta}_{\text{GradientDescent}} = -\lambda g$$

• So combine them:

$$\delta_{\text{DampedNewton}} = -(H + \lambda^{-1}I_d)^{-1}g$$
$$= -\lambda(M + \lambda_d)^{-1}g$$

- $\lambda$  small  $\Rightarrow$  conservative gradient step
- $\lambda$  large  $\Rightarrow$  Newton step

#### USE SECOND DERIVATIVES...



# Levenberg-Marquardt

- Just damped Newton with approximate H
- For a special form of f

$$f(x) = \sum_{i} f_i(x)^2$$

- where  $f_i(x)$  are
  - zero-mean
  - small at the optimum

# 1<sup>ST</sup> DERIVATIVES AGAIN



#### Levenberg Marquardt

- Just damped Newton with approximate H
- For a special form of *f*

$$f(x) = \sum_{i} f_{i}(x)^{2} \qquad 2f(x) f'(x)$$

$$\nabla f(x) = \sum_{i} 2f_{i}(x) \nabla f_{i}(x) \qquad 99^{T}$$

$$\nabla \nabla^{T} f(x) = \sum_{i} \sqrt{(f_{i}(x) \cdot \nabla f_{i}(x))} = \sqrt{7} \int f_{i}(x) + f(x) \nabla f_{i}(x)$$

#### **BACK TO FIRST DERIVATIVES**

**Microsoft** 

Levenberg Marquardt

- Just damped Newton with approximate H
- For a special form of f

$$\begin{split} f(x) &= \sum_{i} f_{i}(x)^{2} \\ \nabla f(x) &= \sum_{i} 2f_{i}(x) \nabla f_{i}(x) \\ \nabla \nabla^{\top} f(x) &= 2 \sum_{i} \mathcal{F}_{i}(x) \nabla \nabla^{\top} f_{i}(x) + \nabla f_{i}(x) \nabla^{\top} f_{i}(x) \end{split}$$

**BACK TO FIRST DERIVATIVES** 





#### CONCLUSION: YMMV

**Microsoft** 

# CONCLUSION: YMMV





For 
$$k=1: 500$$
 f:  $\mathbb{R}^{n} \rightarrow \mathbb{R}$   
 $X_{0} = "randn"(n, 1);$   
 $X^{*} = fmin(f, X_{0});$   
 $E(k) = f(X^{*})$   
end  
plot(sort(E), 1: 500)











FACE



#### Truth

- On many problems, alternation is just fine
  - Indeed always start with a couple of alternation steps
- Computing 2<sup>nd</sup> derivatives is a pain
  - anyone with compiler/parser experience please contact me.

#### BuT

- Just alternation is fine
  - Unless you're willing to problem-select
- Alternation has a convergence guarantee
- Inverting the Hessian is always O(n<sup>3</sup>)

There is a universal optimizer

# TRUTH AND $BuT[\exists u \in \{H, I, L, S, T, U\}^6]$

