Name: Hongda Li
AMATH 516 HW4 2021 FALL

## 4.2

This is the Fenchel Rockafella Primal Dual Problems:

$$(P) \inf_{x \in \mathbf{E}} \{h(\mathcal{A}x) + g(x)\} \qquad (4.2.1)$$

$$(D) \sup_{y \in \mathbf{Y}} \{-h^\star - g^\star(-\mathcal{A}^*y)\}$$

Where we use $\star$ for the Fenchel Conjugate of the function and $*$ for adjoint operator.
Sometimes we need to take the subdifferential/gradient/conjugate of an expression involving multiple parameters, we denote it as the following:

$$\nabla[f(x,y)|x] \quad \partial[f(x,y)|x] \quad [f(x,y)|x]^\star$$

Where the bar notation is telling use which variable we are taking the derivative wrt to the parameter.

### 4.2.1

We wish to compute the Fenchel Rockafella's Dual for the following Optimization Problem:

$$\min_x \left\{ \frac{1}{2} \|Ax - b\|_2^2 + \|x\|_1 \right\}$$

To compute the Dual first identify $f, g$:

$$h(Ax) = \frac{1}{2}\|Ax - b\|_2^2 \qquad (4.2.1.1)$$

$$\implies h(x) = \frac{1}{2}\|x - b\|_2^2$$

$$g(x) = \|x\|_1$$

Considering the Fenchel Conjugate of $h(x)$:

$$h^\star(x) = \sup_y \{\langle x, y \rangle - h(y)\} \qquad (4.2.1.2)$$

$$= \sup_y \left\{ \langle x, y \rangle - \frac{1}{2}\|y - b\|_2^2 \right\}$$

$$\partial\left[ \langle x, y \rangle - \frac{1}{2}\|y - b\|_2^2 \,\Big|\, y \right] = \{x - (y - b)\}$$

$$0 \in \{x - (y^+ - b)\}$$

$$\implies y^+ = x + b$$

$$\implies h^\star(x) = \langle x, x + b \rangle - \frac{1}{2}\|x\|_2^2$$

$$= \langle x, x \rangle + \langle x, b \rangle - \frac{1}{2}\|x\|_2^2$$

$$= \|x\|_2^2 + \langle x, b \rangle - \frac{1}{2}\|x\|_2^2$$

$$h^\star(x) = \frac{1}{2}\|x\|_2^2 + \langle x, b \rangle$$

In this case, the function is continuous and it's concave, therefore a minimum eixsts and if I take the derivative and set it zero, I can get the optimal paramter $y^+$ and find the supremum.

Next we compute the Conjugate for $g(x)$:

$$g(x) = \|x\|_1 \tag{4.2.1.3}$$
$$g^\star(x) = \sup_y \langle \langle x, y \rangle - \|y\|_1 \rangle$$
$$= \sup_y \left\{ \sum_{i=1}^n (x_i y_i - |y_i|) \right\}$$
$$= \sum_{i=1}^n \sup_{y_i} \left\{ x_i y_i - |y_i| \right\}$$

The sup can goes in because it's a summing up discrate expression each using $y_i$, and we assume that $y \in \mathbb{R}^n$, and we can say that:

$$\forall 1 \leq i \leq n \tag{4.2.1.4}$$

$$\sup_{y_i} \left\{ x_i y_i - |y_i| \right\} = \max \left\{ \sup_{y_i \geq 0} \left\{ x_i y_i - y_i \right\}, \sup_{y_i < 0} \left\{ x_i y_i + y_i \right\} \right\}$$

$$\sup_{y_i > 0} \left\{ y_i (x_i - 1) \right\} = \begin{cases} \infty & x_i > 1 \\ 0 & x_i \leq 1 \end{cases}$$

$$\sup_{y_i < 0} \left\{ y_i (x_i + 1) \right\} = \begin{cases} 0 & x_i \geq -1 \\ \infty & x_i < -1 \end{cases}$$

$$\implies \forall 1 \leq i \leq n \quad \sup_{y_i} \left\{ x_i y_i - |y_i| \right\} = \begin{cases} 0 & |x_i| \leq 1 \\ \infty & |x_i| > 1 \end{cases}$$

$$\implies \sum_{i=1}^n \sup_{y_i} \left\{ x_i y_i - |y_i| \right\} = \begin{cases} \infty & \|x\|_\infty > 1 \\ 0 & \|x\|_\infty \leq 1 \end{cases}$$

$$\implies g^\star(x) = \delta\{\|x\|_\infty \leq 1\}(x)$$

The conjugate of $g$ is the indicator function of the unit infinity norm ball.

Now we are ready to use the conjugate function to figure out the Fenchel Rockafella's Dual, we use the results from (4.2.1.2) and (4.2.1.4):

$$\sup_y \left\{ -h^\star(y) - g^\star(-\mathcal{A}^* y) \right\} \tag{4.2.1.5}$$

$$\equiv \sup_y \left\{ -\frac{1}{2} \|y\|_2^2 - \langle y, b \rangle - \delta\{\| - A^T y\|_\infty \leq 1\}(x) \right\}$$

$$\equiv \sup_{\|A^T y\|_\infty \leq 1} \left\{ -\frac{1}{2} \|x\|_2^2 - \langle x, b \rangle \right\}$$

## 4.2.2

We wish to look for the Dual for the following problem:

$$\min_{\|x\|_q \leq 1} \|Ax - b\|_p$$

Where we are going to use $\|x\|_m, \|x\|_{\bar{m}}$ to denote the norm and its dual norm. First we identify the $f, g$ as:

$$h(x) = \|x - b\|_p \tag{4.2.2.1}$$
$$g(x) = \delta\{\|x\|_q \leq 1\}(x)$$

2

Before we start I am going to introduce a convenient new trick and I am going to prove it:

**Claim**: "Conjugate of a Norm is the indicator function for the unit norm ball of the Dual Norm."

**Proof**:

Set $f(x) = \|x\|$ and we use $\|x\|_*$ to denote the dual norm, and we wish to prove that $f^\star(x) = \delta\{\|x\|_* \leq 1\}(x)$. Recall that the definition of dual norm is:

$$\|y\|_* = \sup_{\|x\| \leq 1} \langle x, y \rangle$$

From the definition of the conjugate of a function we have: $f^\star(x) = \sup_y\{\langle x, y \rangle - f(y)\}$, and let's discuss it by cases:

Suppose that $\|y\|_* \leq 1$ then:

$$1 \geq \|y\|_* \geq \sup_{\|x\| \leq 1} \{\langle x, y \rangle\} \tag{4.2.2.2}$$

$$\forall \|x\| \leq 1 : \langle x, y \rangle \leq 1$$

$$\text{let: } t > 0$$

$$\langle tx, y \rangle \leq t\langle x, y \rangle \leq t = t\|x\| = \|tx\|$$

$$\implies \forall\, tx\, , \|x\| \leq 1 : \langle tx, y \rangle \leq \|tx\|$$

$$\implies \forall x\, \langle x, y \rangle \leq \|x\|$$

$$\implies \sup_x \{\langle y, x \rangle - \|x\|\} = 0$$

Because the inner produce is forever smaller than the norm of $x$, therefore we have to set the value of $x = 0$ to attain the sup of the expression, which is zero.

Next we consider the case where $\|y\|_* > 1$, and we can say that:

$$\|y\|_* > 1 \tag{4.2.2.3}$$

$$\exists x : \|x\| \leq 1 \wedge \langle x, y \rangle > 1$$

$$\implies \sup_y \{\langle x, y \rangle - \|x\|\} \geq \langle x, y \rangle - \|x\| > 0$$

$$\text{let } x = tx, t \geq 0$$

$$t\langle x, y \rangle - t\|x\| > 0$$

$$t(\langle x, y \rangle - \|x\|) > 0$$

$$\implies \sup_y \{\langle x, y \rangle - \|x\|\} = \infty$$

Therefore the conjugate of the function $f$ is: $f^\star(x) = \delta\{\|x\|_* \leq 1\}(x)$. It's the indicator function of the dual unit norm ball.

Using this pieces of information, we can start figuring out the conjugate. From table 3.4 and the fact that the function $f, h$ are both norm of something, hence convex, therefore the double conjugate of the function is the funciton itself. Therefore the table can be applied in both forward and backwards direction. In our case, we are applying the seceond entry of table 3.4.

$$[f(x + b)\,|\,x]^\star = f^\star(y) - \langle b, y \rangle \tag{4.2.2.4}$$

In our case we have:

$$[\,\|x - b\|_p\,|\,x]^\star (y) = \delta\{\|x\|_{\bar{p}} \leq 1\}(y) - \langle b, y \rangle \tag{4.2.2.5}$$

By a similar token the conjugate of the $g$ function would be:

$$[\delta\{\|x\|_{q \leq 1}(x)\}]^\star (y) = \|y\|_{\bar{q}} \tag{4.2.2.6}$$

3

Results are out:
$$h(x) = \|x - b\|_p \quad g(x) = \delta\{\|q\| \leq 1\}(x)$$

And using the Rockafella's Duality we have:
$$\sup_{\|y\|_{\bar{p}} \leq 1} \left\{ \langle b, y \rangle - \|A^T y\|_{\bar{q}} \right\} \tag{4.2.2.7}$$

The negative sign can get ignore because $g^\star$ is a norm.

### 4.2.3

The primal is:
$$(P) : \min_x \left\{ \langle c, x \rangle : Ax = b, x \in K \right\}$$

Where $K$ is a cone.
Identifies the $g, h$ function in this case it's:
$$h(Ax) = \delta\{Ax = b\} \tag{4.2.3.1}$$
$$\implies h(x) = \delta\{x = b\}$$
$$g(x) = \langle c, x \rangle + \delta\{x \in K\}$$

We can immdiately identifies the conjugate of the indicator function $h$ because it's part of the Exercise 3.24 from the last homework. And therefore we have: $h^\star(x) = \langle b, x \rangle$.
To figure out the Conjugate of $h$:
$$g^\star(x) = \sup\{\langle x, y \rangle - g(y)\} \tag{4.2.3.2}$$
$$= \sup_y \{\langle x, y \rangle - \langle c, y \rangle - \delta\{y \in K\}(y)\}$$
$$= \sup_y \{\langle y, x - c \rangle - \delta\{y \in K\}(y)\}$$
$$= \sup_{y \in K} \{\langle y, x - c \rangle\}$$
$$= \delta_K^\star(x - c)$$

Using page 60 of the textbook, and the fact that $K$ is a cone, we have: $g^\star(x) = \delta_K^\star(x - c) = \delta_{K\circ}(x - c)$.
The Fenchel Dual is:
$$\sup_y \left\{ -h^\star(y) - g^\star(-A^*y) \right\} = \sup_y \left\{ -\langle b, y \rangle - \delta_{K\circ}(-A^T y - c) \right\} \tag{4.2.3.3}$$
$$\equiv \sup_y \left\{ \langle b, y \rangle - \delta_{K\circ}(A^T y - c) \right\}$$

### 4.2.4

The Primal Problem is:
$$(P) : \min_x \left\{ \frac{1}{2} \langle x, Qx \rangle + \langle c, x \rangle : Ax \leq b \right\}$$

Where the matrix $Q$ is Positive Definite. Identifying the $g, h$ we have:
$$h(Ax) = \delta\{Ax \geq b\}(x) \tag{4.2.4.1}$$
$$g(x) = \frac{1}{2}\|x\|_Q^2 + \langle x, c \rangle$$
$$\implies h(x) = \delta\{x \geq b\}(x)$$
$$= \delta\{x - b \geq 0\}(x)$$

4

Notice that $x - b \geq 0 \iff x - b \in \mathbb{R}_+^n$ and $\mathbb{R}_+^n$ is a cone and it's polar cone is $\mathbb{R}_-^n$. Therefore, we use table 3.4:

$$h(x) = \delta_{\mathbb{R}_+^n}(x - b) \tag{4.2.4.2}$$

$$\implies h(x + b) = \delta_{\mathbb{R}_+^n}(x)$$

$$\implies h^\star(x) = \delta_{\mathbb{R}_-^n}(x) - \langle b, x \rangle$$

The conjugate of $g$ is given by:

$$g^\star(y) = \left[ \frac{1}{2}\|x\|_Q^2 + \langle c, x \rangle \,\Big|\, x \right]^\star (y) \tag{4.2.4.3}$$

$$= \left[ \frac{1}{2}\|x\|_Q^2 \,\Big|\, x \right]^\star (y - c)$$

Let's figure out the conjugate of the one half of the Energy Norm squared:

$$\left[ \frac{1}{2}\|A^{1/2}x\|_2^2 \right]^\star (y) = \sup_x \left\{ \langle y, x \rangle - \frac{1}{2}\|A^{1/2}x\|_2^2 \right\} \tag{4.2.4.4}$$

Set Gradient to zero to find the sup

$$\nabla \left[ \langle y, x \rangle - \frac{1}{2}\|x\|_A^2 \,\Big|\, x \right] = 0$$

$$\implies y - Ax^+ = 0$$

$$Ax^+ = y$$

$$x^+ = A^{-1}y$$

$$\implies \sup_x \left\{ \langle y, x \rangle - \frac{1}{2}\|A^{1/2}x\|_2^2 \right\} = \langle y, A^{-1}y \rangle - \frac{1}{2}\|A^{-1/2}y\|_2^2$$

$$= \|y\|_{A^{-1}}^2 - \frac{1}{2}\|y\|_{A^{-1}}^2$$

$$= \frac{1}{2}\|y\|_{A^{-1}}^2$$

Then we can conclude together with the results from (4.2.4.3) to get:

$$g^\star(y) = \frac{1}{2}\|y - c\|_{Q^{-1}}^2 \tag{4.2.4.5}$$

Therefore, the Dual is:

$$\max_y \left\{ -h^\star(y) - g^\star(-\mathcal{A}^* y) \right\} = \max_y \left\{ \delta_{\mathbb{R}_-^n}(y) - \langle b, y \rangle - \frac{1}{2}\| - A^T y - c\|_{Q^{-1}}^2 \right\} \tag{4.2.4.6}$$

$$\equiv \max_{y \geq 0} \left\{ \langle b, y \rangle - \frac{1}{2}\|A^T y - c\|_{Q^{-1}}^2 \right\}$$

## 5.14

For this part of the HW I implemented the Smooth Gradient Descend and the Accelerated Gradient on Julia the programming language and made the plot of for the values of the objective functions for the first 100 iterations of both method. (Running the code requires installing UnicodePlots.jl and Plots.jl)
This is how I got the gradient for the function for the implementations.
Let $X \in \mathbb{R}^{m \times n}$, let $y \in \mathbb{R}^m$, $\theta \in \mathbb{R}^n$. And let lower case $x_i$ denotes the $i$ th of the matrix $X$ The Gradient

of the loss function is computed as:

$$\nabla_\theta \left[ \sum_{i=1}^{m} \ln(1 + \exp\left(\langle\theta, -y_i x_i\rangle\right))) + \frac{\lambda}{2}\|\theta\|_2^2 \right] \tag{5.14.1}$$

$$= \left( \sum_{i=1}^{m} \nabla_\theta \left[\ln\left(1 + \exp\langle\theta, -y_i x_i\rangle\right)\right] \right) + \lambda\theta$$

$$= \left( \sum_{i=1}^{m} \frac{\exp(\theta, -yx_i)}{1 + \exp(\theta, -yx_i)}\nabla_\theta[\langle\theta, -y_i x_i\rangle] \right) + \lambda\theta$$

$$= \left( \sum_{i=1}^{m} \frac{-\exp(\theta, -yx_i)}{1 + \exp(\theta, -yx_i)}yx_i \right) + \lambda\theta$$

The code implementations is here:

```
1   using LinearAlgebra
2   import UnicodePlots
3   import Plots
4   import Logging
5
6   # -----------------------------------------------------------------------------
7   mutable struct LogisticRegressionModel
8       X::Matrix{Float64}      # The data
9       y::Vector{Float64}      # the labels
10      theta::Vector{Float64}  # The parameters
11      lambda::Float64         # a regularization parameters
12
13      function LogisticRegressionModel(X, y, theta; lambda=1)
14          this = new(X, y, theta, lambda)
15          return this
16      end
17
18  end
19
20  function Base.copy(this::LogisticRegressionModel)
21      return LogisticRegressionModel(
22          copy(this.X),
23          copy(this.y),
24          copy(this.theta),
25          lambda = copy(this.lambda)
26          )
27  end
28
29  function ErrorOf(this::LogisticRegressionModel)
30      theta = this.theta
31      X = this.X
32      y = this.y
33      lambda = this.lambda
34
35      Error = - y.*(X*theta)  # vector
36      Error = 1 .+ exp.(Error)  # vector
37      Error .= log.(Error)   # vector
38      Error = sum(Error)
39      Error += (lambda/2)*dot(theta, theta)
```

6

```julia
40          return Error
41     end
42
43     function GradientOf(this::LogisticRegressionModel)
44          theta = this.theta
45          X = this.X
46          y = this.y
47          lambda = this.lambda
48
49          Expr1 = (x) -> exp(x)/(1 + exp(x))
50          Summation = zeros(size(theta))
51          for Index = 1:size(X, 1)
52               x = X[Index, :]
53               Expr2 = - y[Index]*x   # vector
54               Expr3 = dot(theta, Expr2)   # scalar
55               Summation += Expr1(Expr3) * Expr2   # Vector
56          end
57          # ----------------------------------------------
58          # Summation = exp.(-y.*X*theta)
59          # Summation = Summation./(Summation .+ 1)
60          # Summation = X'*(-y.*Summation)
61          return Summation + lambda*theta
62     end
63
64     function ParametersOf(this::LogisticRegressionModel)
65          return this.theta
66     end
67
68
69
70     # ----------------------------------------------------------------------------
71     mutable struct GradientDescend
72          beta::Float64
73          objective::LogisticRegressionModel
74          objective_vals::Vector{Float64}
75
76          function GradientDescend(objective::LogisticRegressionModel, beta::Float64)
77               this = new(beta, objective, Vector{Float64}())
78               push!(this.objective_vals, ErrorOf(this.objective))   # initial error.
79               return this
80          end
81     end
82
83
84     function GradientUpdate(this::GradientDescend)
85          DeltaTheta = GradientOf(this.objective)
86          Params = ParametersOf(this.objective)
87          Params .-= DeltaTheta * (1/this.beta)
88          push!(this.objective_vals, ErrorOf(this.objective))
89
90          return norm(DeltaTheta)
91     end
92
93     # ----------------------------------------------------------------------------
```

```julia
94   mutable struct AcceleratedGradientDescend
95       beta::Float64
96       objective::LogisticRegressionModel
97       objective_vals::Vector{Float64}
98       t::Int64                        # Iteration number.
99       a::Dict{Int64, Float64}
100      x::Dict{Int64, Vector{Float64}}
101
102      function AcceleratedGradientDescend(objective::LogisticRegressionModel, beta::Float64)
103          this = new(
104              beta,
105              objective,
106              Vector{Float64}(),
107              0,
108              Dict{Int64, Float64}(),
109              Dict{Int64, Vector{Float64}}()
110          )
111          push!(this.objective_vals, ErrorOf(this.objective))  # initial error.
112          a = this.a
113          x = this.x
114          a[-1] = a[0] = 1
115          x[-1] = x[0] = this.objective.theta
116          return this
117      end
118  end
119
120  function GradientUpdate(this::AcceleratedGradientDescend)
121      t = this.t
122      x = this.x
123      a = this.a
124      beta = this.beta
125      u = x[t] + a[t]*(a[t - 1]^(-1) - 1)*(x[t] - x[t - 1])
126      Gradient = GradientOf(this.objective)
127      x[t + 1] = u - (1/beta)*Gradient
128      this.objective.theta .= x[t + 1]
129      a[t + 1] = (sqrt(a[t]^4 + 4a[t]^2) - a[t]^2)/2
130      this.t += 1
131      push!(this.objective_vals, ErrorOf(this.objective))
132      return norm(Gradient)
133  end
134
135
136
137  function Run()
138      function Bernoulli(p)
139          if rand() < p
140              return 1
141          end
142          return -1
143      end
144      m, n = 100, 50
145      beta = 100.0
146      theta = ones(n)
147      X = randn(m, n)
```
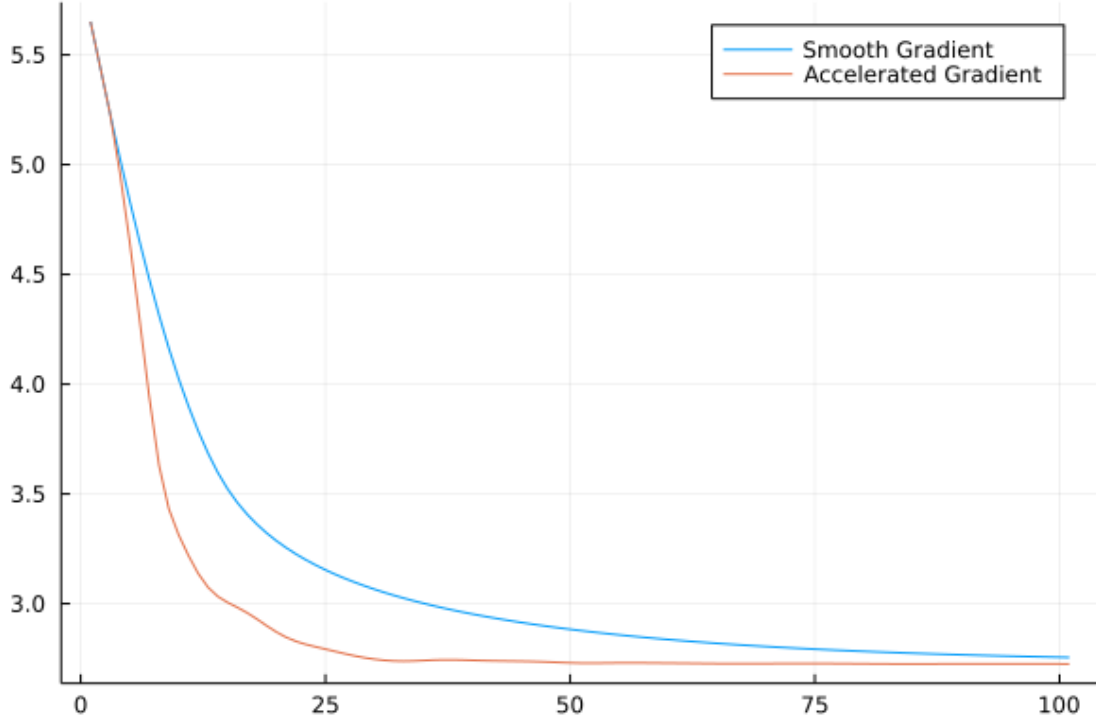
```
148     z = X*theta
149     p = 1 ./ (1 .+ exp.(-z))
150     theta0 = randn(n)
151     y = Bernoulli.(p)
152     Objective = LogisticRegressionModel(X, y, theta0)
153     Optim = GradientDescend(Objective, beta)
154     OptimAcc = AcceleratedGradientDescend(copy(Objective), beta)
155     println("running...")
156     for _ in 1:100
157         println(GradientUpdate(Optim))
158         GradientUpdate(OptimAcc)
159     end
160     Logging.@info "Objective Values (Smooth Gradient Descend):"
161     FxnVals = Optim.objective_vals
162     FxnValsAcc = OptimAcc.objective_vals
163     display(FxnVals)
164     Logging.@info "Objective Values (Acc Gradient Descend)"
165     display(FxnValsAcc)
166
167     # Plot this out -----------------------------------------------------------
168     UPlot = UnicodePlots.lineplot(1:length(FxnVals), FxnVals)
169     println(UnicodePlots.lineplot!(UPlot, 1:length(FxnValsAcc), FxnValsAcc))
170     Plots.plot(collect(1:length(FxnVals)), log.(FxnVals), label="Smooth Gradient")
171     Plots.plot!(collect(1:length(FxnValsAcc)), log.(FxnValsAcc), label="Accelerated Gradien
172
173     Plots.savefig("objectiveVals.png")
174
175     return Optim
176 end
177
178 Optim = Run()
179
180
```

The plot produced by the algorithms is here(Next page):

Hi Chatherine, if you need the file for the code and the Julia Project Environment, please contact me I am happy to provide.

## 5.15

### 5.15.1

We assume that the $f(x)$ is convex differentiable with minimizer $x^\star$. And we wish to show the following:

$$\text{Gradient Update: } x_{k+1} = x_k + \lambda_k \nabla f(x_k) \tag{5.15.1.1}$$

$$\text{Wish to prove: } \frac{1}{2}\|x_{k+1} - x_*\|^2 \leq \frac{1}{2}\|x - x_*\|^2 - \gamma_k(f(x_k) - f(x_*)) + \frac{\gamma_k}{2}\|\nabla f(x_k)\|$$

To start consider that:

$$\|(x_{k+1} - x_k) - (x_k - x_*)\|^2 = \|x_{k+1}\|^2 + \|x_k - x_*\|^2 + 2\langle x_{k+1} - x_k, x_k - x_*\rangle \tag{5.15.1.2}$$
$$= \|x - x_*\|^2 + \|\gamma_k \nabla f(x_k)\|^2 + 2\langle \gamma_k \nabla f(x_k), x_k - x_*\rangle$$

Divides the whole expression by 2 and focuses on the last 2 terms of the expression:

$$\frac{1}{2}\|\gamma_k \nabla f(x_k)\|^2 + \langle \gamma_k \nabla f(x_k), x_k - x_*\rangle \tag{5.15.1.3}$$
$$= \frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 + \gamma_k\langle \nabla f(x_k), x_k - x_*\rangle$$

Using Excercise 4.12 #4 we have:

$$\langle \nabla f(x_k) - \underbrace{\nabla f(x_*)}_{=0}, x_* - x_k\rangle \geq 0 \tag{5.15.1.4}$$
$$f(x_*) - f(x_k) \geq \langle \nabla f(x_k), x_k - x_*\rangle$$

10

Apply (5.15.1.4) to (5.15.1.3) we have:

$$\frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 + \gamma_k\langle\nabla f(x_k), x_k - x_*\rangle \le \frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 + \gamma_k(f(x_*) - f(x_k)) \qquad (5.15.1.5)$$

Go back to (5.15.1.2) and apply that previous bound we have:

$$
\begin{aligned}
\frac{1}{2}\|(x_{k+1} - x_k) - (x_k - x_*)\|^2 &= \frac{1}{2}\|x - x_*\|^2 + \frac{1}{2}\|\gamma_k \nabla f(x_k)\|^2 + \langle\gamma_k \nabla f(x_k), x_k - x_*\rangle \qquad (5.15.1.6)\\
&\le \frac{1}{2}\|x - x_*\|^2 + \frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 + \gamma_k(f(x_*) - f(x_k))\\
&= \frac{1}{2}\|x - x_*\|^2 + \frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 - \gamma_k(f(x_k) - f(x_*))
\end{aligned}
$$

Which is exactly what we want to show.

### 5.15.2

We wish to show that that sequence is minimized by choosing $\gamma_k = (f(x_k) - f^*)/\|\nabla f(x_k)\|^2$, where $f^*$ is the minimal value for the function. This is what we wish to prove:

$$\|x_{k+1} - x^*\|^2 \le \|x_k - x^*\|^2 - \left(\frac{f(x_k) - f^*}{\|\nabla f(x_k)\|}\right)^2$$

And then proving the the sequence is bounded by:

$$
\begin{aligned}
\frac{1}{2}\|x_{k+1} - x_*\|^2 &\le \frac{1}{2}\|x_k - x_*\|^2 - \gamma_k(f(x_k) - f^*) + \frac{\gamma_k^2}{2}\|\nabla f(x_k)\|^2 \qquad (5.15.2.1)\\
&= \frac{1}{2}\|x_k - x_*\|^2 - \gamma_k\left(f(x_k) - f^* - \frac{\gamma_k}{2}\|\nabla f(x_k)\|^2\right)
\end{aligned}
$$

$$\text{Let: } \partial_{\gamma_k}\left[\gamma_k(f(x_k) - f^*) + \frac{\gamma^2}{2}\|\nabla f(x_k)\|^2\right] = 0$$

$$-(f(x_k) - f^*) + \gamma_k^+\|\nabla f(x_k)\|^2 = 0$$

$$\gamma_k^+ = \frac{f(x_k) - f^*}{\|\nabla f(x_k)\|^2}$$

Consider the RHS of the first expression in the above block of statements, and substitute the optimal stepside of the expression we have:

$$
\begin{aligned}
\|x_{k+1} - x_*\|^2 &\le \|x_k - x_*\|^2 - 2\gamma_k^+(f(x_k) - f^*) + (\gamma_k^+)^2\|\nabla f(x_k)\|^2 \qquad (5.15.2.2)\\
&= \|x_k - x_*\|^2 - \gamma_k^+(2(f(x_k) - f^*) - \gamma^+\|\nabla f(x_k)\|^2)\\
&= \|x_k - x_*\|^2 - \gamma_k^+(2(f(x_k) - f^*) - (f(x_k) - f^*))\\
&= \|x_k - x_*\| - \gamma_k^+(f(x_k) - f^*)\\
&= \|x_k - x_*\| + \left(\frac{f(x_k) - f^*}{\|\nabla f(x_k)\|}\right)^2
\end{aligned}
$$

Which is the statement that we wish to prove.

### 5.15.3

In this part we wish to prove that :

$$f\left(\frac{1}{k}\sum_{i=0}^{k-1} x_i\right) - f^* \le \frac{2\beta\|x_0 - x_*\|^2}{k}$$

11

Inaddition when the function is $\alpha$ strongly convex, wehave the additional results of:

$$\|x_{k+1} - x_*\|^2 \leq \left(1 - \frac{\alpha}{4\beta}\right)\|x_k - x_*\|^2$$

**Proof**: To make our life easier, we consider using the following notations for the errors and the optimality gap:

$$E_k = f(x_k) - f^* \tag{5.15.3.1}$$
$$e_k = x_k - x^*$$

Starting with exercise 3.12, the differential characterization for beta smooth function, we have:

$$f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2\beta}\|\nabla f(x) - \nabla f(y)\|^2 \leq f(y) \tag{5.15.3.2}$$

$$\text{consider: } y = x_k \quad x = x^*$$

$$f^* + \frac{1}{2\beta}\|\nabla f(x_k)\|^2 \leq f(x_k)$$

$$\|\nabla f(x_k)\|^2 \leq 2\beta(f(x_k) - f^*) = 2\beta E_k$$

$$\implies \frac{1}{2\beta} \leq \frac{E_k}{\|\nabla f(x_k)\|^2}$$

$$\implies \frac{E_k}{2\beta} = \frac{E_k^2}{\|\nabla f(x_k)\|^2}$$

Then we start with the results from what we proved in (5.15.2.2), And written it as:

$$\|e_{k+1}\|^2 \leq \|e_k\|^2 - \left(\frac{E_k}{\|\nabla f(x_k)\|}\right)^2 \tag{5.15.3.3}$$

$$\left(\frac{E_k}{\|\nabla f(x_k)\|}\right)^2 \leq \|e_k\|^2 - \|e_{k+1}\|^2$$

$$\underset{(5.15.3.2)}{\implies} \frac{E_k}{2\beta} \leq \|e_k\|^2 - \|e_{k+1}\|^2$$

Consider summing up the terms in the above inequalities for $0 \leq i \leq k - 1$:

$$\sum_{i=0}^{k-1} \frac{E_i}{2\beta} \leq \sum_{i=1}^{k-1} \|e_i\|^2 - \|e_{i+1}\|^2 \tag{5.15.3.4}$$

$$\frac{1}{k}\sum_{i=1}^{k-1} E_i \leq \frac{2\beta}{k}\sum_{i=1}^{k-1} \|e_i\|^2 - \|e_{i+1}\|^2$$

$$\frac{1}{k}\sum_{i=0}^{k-1} E_i \leq \frac{2\beta}{k}\left(\|e_0\|^2 - \|e_k\|^2\right)$$

$$\underset{\text{By Convexity}}{\implies} f\left(\frac{1}{k}\sum_{i=0}^{k-1} E_i\right) - f^* \leq \frac{2\beta}{k}\left(\|e_0\|^2 - \|e_k\|^2\right)$$

$$f\left(\frac{1}{k}\sum_{i=0}^{k-1} E_i\right) - f^* \leq \frac{2\beta}{k}\|e_0\|^2$$

$$f\left(\frac{1}{k}\sum_{i=0}^{k-1} x_i\right) - f^* \leq \frac{2\beta\|x_0 - x_*\|^2}{k}$$

I used telescoping from the second line to the first line in the above block. Just expand the sume and the adjacent terms of the differences will get canceled out.

12

From the second last line to the last line, we can remove the $-\|e_k\|^2$ becase it's larger than zero and removing it will make the quantity larger.

Now let's consider the statement from 3.57, the differential characterization of $\alpha$ strongly convex function:

$$f(x_k) - f^* \geq \frac{\alpha}{2}\|x_k - x^*\|^2 \tag{5.15.3.5}$$

$$\|\nabla f(x_k)\| = \|\nabla f(x_k) - \underbrace{\nabla f(x_*)}_{=0}\| \leq \beta\|x_k - x^*\| \tag{5.15.3.6}$$

Now we may reconsider the results from (5.15.3.4), which is giving us:

$$\|x_{k+1} - x_*\|^2 \leq \|x_k - x_*\|^2 - \left(\frac{f(x-k) - f^*}{\|\nabla f(x_k)\|}\right) \tag{5.15.3.7}$$

$$\|x_{k+1} - x_*\|^2 \leq \|x_k - x^*\|^2 - \left(\frac{\frac{\alpha}{2}\|x_k - x^*\|^2}{\beta\|x_k - x_*\|^2}\right)$$

$$= \left(1 - \frac{\alpha^2}{4\beta^2}\right)\|x_k - x_*\|^2$$

From the first expression to the second one, I replaces the numerator with a smaller quantity, the denominator with a larger quantity, reducing the size of the fraction as whole, and hence reducing the quantity being subtracted smaller, placing an upper bound on the RHS of the first inequality.

That last expression is what we wish to prove for the problem.

## 6.11

We were given the accelerated Gradient Descend method which has an interesting sequence that weights on the previous velocity of the descend steps.

### 6.11.1

We wish to show that this relation for the sequence holds:

$$\frac{1 - a_{t+1}}{a_{t+1}^2} = \frac{1}{a_t^2} \quad \forall t \geq 0 \tag{6.11.1.1}$$

**Proof**:

$$a_{t+1} = \frac{1}{2}\left(\sqrt{a_t^4 + 4a_t^2} - a_t^2\right) \tag{6.11.1.2}$$

$$2a_{t+1} = \sqrt{a_t^4 + 4a_t^2} - a_t^2$$

$$2a_{t+1} + a_t = \sqrt{a_t^4 + 4a_t^2}$$

$$(2a_{t+1} + a_t^2)^2 = a_t^4 + 4a_t^2$$

$$4a_{t+1}^2 + a_t^4 + 4a_t^2 a_{t+1} = a_t^4 + 4a_t^2$$

$$4a_{t+1}^2 + 4a_{t+1}a_t^2 - 4a_t^2 = 0$$

$$a_{t+1}^2 + a_{t+1}a_t^2 - a_t^2 = 0$$

$$a_{t+1}^2 + a_t^2(a_{t+1} - 1) = 0$$

$$\frac{1}{a_t^2} + \frac{a_{t+1} - 1}{a_{t+1}^2} = \frac{0}{a_t^2 a_{t+1}^2}$$

$$\frac{a_{t+1} - 1}{a_{t+1}^2} = \frac{1}{a_t^2} \quad \forall t \geq 0$$

13

### 6.11.2

We wish to prove that $\sum_{i=0}^{t} \frac{1}{a_i} = \frac{1}{a_t^2}$, and then the sequence $a_k$ is upper bounded by the sequence $2/(t+2)$. Inductively assume that the statement holds for $n$:

$$\sum_{i=0}^{t} \frac{1}{a_i} = \frac{1}{a_t^2} \quad \forall t \leq n \tag{6.11.2.1}$$

$$\frac{1 - a_{t+1}}{a_{t+1}^2} = \frac{1}{a_t^2} \tag{6.11.2.2}$$

$$\frac{1}{a_{t+1}^2} - \frac{1}{a_{t+1}} = \frac{1}{a_t^2}$$

$$\frac{1}{a_{t+1}^2} = \frac{1}{a_t^2} + \frac{1}{a_{t+1}}$$

$$\frac{1}{a_{t+1}^2} = \sum_{i=0}^{t} a_i + \frac{1}{a_{t+1}}$$

$$\frac{1}{a_{t+1}^2} = \sum_{i=0}^{t+1} \frac{1}{a_i}$$

The statement paramaterized by $n+1$ holds true. The base case it's not hard to verify, just compute:

$$\frac{1}{a_1} + \frac{1}{a_0} = \frac{1}{\frac{1}{2}(\sqrt{5} - 1)} + 1 \tag{6.11.2.3}$$

$$\frac{1}{a_1^2} = \frac{1}{\frac{1}{4}(\sqrt{5} - 1)^2}$$

Which is actually proved by the end of (6.11.1.2).
Next, we wish to prove the upper bound for the sequence.

This problem is too hard and I can't prove it. I accessed the solutions to multiple peers and they turned out to contain their own algebraic errors. I don't think the upper bound is proven via induction. I do know that the sequence approaches zero monotone. And judging by the amount of efforts I put into this problem, I don't think it's worth it.

Inductively, I used a lot of methods and the fact that the sequence is monotone decreasing.

I gave up after discovering something huge. I try to insert epsilons in inequalities trying to look for additional assumptions can make it better, however, it's futile because all epsilon has to be less than zero, proving that all inductive proof is likely to fail.

Therefore, the proof for the question is not via induction, and it's very likely that this is true.

I computationally verify that the upper bound is extremely tight, and the difference between the sequence $2/(k+2)$ and $a_k$ is a power sequence of power $-2$. The bound is very tight.

Hi Catherine, please check other's algebra on this problem carefully, the proof you went through on Thursday from a student has an algebraic mistake. If you find any legit proof for this problem, please share it via Canvas message to me, I want to know how to handle such a tight bound for this sequence.

Here is the code I used for the numerical investigation of the problem. You will need it to intuitively verify the correct solutions. The code is in matlab.

```
1   close all;
2   a = 1;
3   Points = 0:300;
4   for Index = Points(2:end)
5       p = a(end); % Previous
```

```matlab
 6        a(end + 1) = (sqrt(p^4 + 4*p^2) - p^2)/2;
 7    end
 8    figure;
 9    plot(a, '-o', "linewidth", 1); hold on
10    UpperBound = 2./(Points + 2);
11    LowerBound = 2./(Points + 3);
12    plot(UpperBound, "linewidth", 1);
13    plot(LowerBound, "linewidth", 1);
14    legend(["sequence", "upper bound", "Lower Bound"])
15
16    figure;
17    DiffBand = UpperBound - LowerBound;
18    loglog(DiffBand);
19    Coeff = polyfit(log(Points(2:end)), log(DiffBand(2:end)), 1);
20    disp("Log log Coefficient:")
21    disp(num2str(Coeff(1)))
```