

Q1 "Just make a nicer plot"

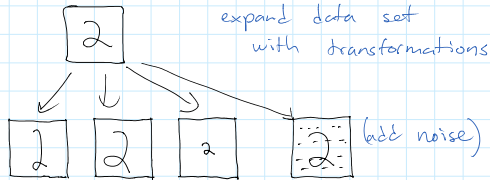
Q2 should be easy to obtain 0.25 - 0.15

1. more hidden units
2. play with step size, momentum
3. creating more training examples may help the most

Image toolbox:

- 'im resize
- 'im rotate
- 'im translate

```
I = reshape(X(i,:)[16,16])
imageSc(I)
```



Q3 $\mu I \leq \nabla^2 f(x)$ $\nabla_{ii}^2 f(x) \leq L$

an "ordering" on the matrices

$$y^T \nabla^2 f(x) y \geq y^T (\mu I) y$$

$$\downarrow$$

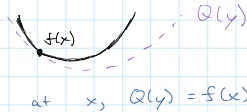
$$y^T \nabla^2 f(x) y \geq \mu y^T y$$



diagonals less than L

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)$$

here exists some z that this holds true



set y to (y-x)

$$(y-x)^T \nabla^2 f(z) (y-x) \geq \mu (y-x)^T (y-x) \quad \frac{\mu}{2} \|y-x\|^2$$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} (y-x)^T (y-x)$$

$$\left\{ a^T y - c + \frac{\mu}{2} y^T y - \mu y^T x + c \right\}$$

take derivative, set to zero

minimize wrt y:

$$\nabla Q(y) = 0 + \nabla f(x) + \mu (y-x)$$

solve: $y^* = x - \left(\frac{1}{\mu}\right) \nabla f(x)$

$$f(y^*) \geq f(x) + \underbrace{\nabla f(x)^T \left(x - \frac{1}{\mu} \nabla f(x)\right) - x}_{-\frac{1}{\mu} \nabla f(x)^T \nabla f(x)} + \frac{\mu}{2} \left\| \left(x - \frac{1}{\mu} \nabla f(x)\right) - x \right\|^2$$

$$-\frac{1}{\mu} \nabla f(x)^T \nabla f(x) \Rightarrow -\frac{1}{\mu} \|\nabla f(x)\|^2$$

$$f(x^*) \geq f(x) - \frac{1}{\mu} \|\nabla f(x)\|^2 + \frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$f(x^*) - f(x) \geq -\frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$-2\mu [f(x) - f(x^*)] \geq \|\nabla f(x)\|^2$$

$$f(y) \geq g(y)$$

$$f(y^*) \geq g(y^*)$$

↑ minimizes f ↑ minimizes g

$$g(y^*) \geq g(y^*)$$

don't need proof

$$y = x^{t+1} = x^t - \frac{1}{L} \nabla_{i_t} f(x^t) e_{i_t}$$

pick coordinate i_t to update

$$= x^t - \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \frac{1}{L} \nabla_{i_t} f(x^t) \\ \vdots \\ 0 \end{bmatrix}$$

$$e_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i$$

$y = x$, except for coordinate i assume only differ in one coordinate

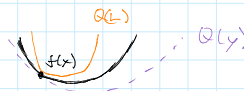
$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)$$

$$f(y) = f(x) + \nabla_i f(x) (y-x)_i + \frac{1}{2} (y-x)_i^2 \nabla_{ii} f(z)$$

$$\nabla_{ii} f(x) \leq L$$

all multiplied by 0
except at coordinate i

$$f(y) \leq f(x) + \nabla_i f(x) (y-x)_i + \frac{1}{2} (y-x)_i^2$$



set $x = x^t$
 $y = x^{t+1}$

$$f(x^{t+1}) \leq f(x^t) + \nabla_i f(x^t) (x^{t+1} - x^t)_i + \frac{L}{2} (x^{t+1} - x^t)_i^2$$

$$y = x^{t+1} = x^t - \frac{1}{L} \nabla_{i_t} f(x^t) e_{i_t}$$

$$(x^{t+1} - x^t) = -\frac{1}{L} \nabla_{i_t} f(x^t) e_{i_t}$$

$$f(x^{t+1}) \leq f(x^t) + \nabla_i f(x^t) \left(-\frac{1}{L} \nabla_{i_t} f(x^t) e_{i_t}\right)_i + \frac{L}{2} \left\| \frac{1}{L} \nabla_{i_t} f(x^t) e_{i_t} \right\|^2$$

$$\Rightarrow f(x^t) - \frac{1}{L} (\nabla_i f(x^t))^2 + \frac{1}{2L} (\nabla_i f(x^t))^2$$

$$\Rightarrow f(x^t) - \frac{1}{2L} (\nabla_i f(x^t))^2$$

\Rightarrow this is for all methods (all 3 parts of question)

Coordinate Descent:

upper bound: $f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} (\nabla_i f(x^t))^2$

lower bound: $f(x^*) \geq f(x^t) - \frac{1}{2\mu} \|\nabla f(x^t)\|^2$

Review from class: Gradient Descent:

upper bound $f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2$

lower bound: $f(x^*) \geq f(x^t) - \frac{1}{2\mu} \|\nabla f(x^t)\|^2$
(lower bound is the same)

$$-\frac{1}{2\mu} \|\nabla f(x^t)\|^2 \leq f(x^*) - f(x^t)$$

$$\|\nabla f(x^t)\|^2 \leq -2\mu [f(x^t) - f(x^*)]$$

$$f(x^{t+1}) \leq f(x^t) - \frac{2\mu}{2L} [f(x^t) - f(x^*)]$$

subtract $f(x^*)$ from both sides (a-b > b-c)

$$f(x^{t+1}) - f(x^*) \leq f(x^t) - f(x^*) - \frac{\mu}{L} (f(x^t) - f(x^*))$$

$$\Rightarrow \left(1 - \frac{\mu}{L}\right) [f(x^t) - f(x^*)]$$

$i_t = \operatorname{argmax}_i |\nabla_i f(x^t)|$ biggest improvement by changing
Greedy:

$$(\nabla_{i_t} f(x^t))^2 = \|\nabla f(x^t)\|_\infty^2$$

$$\geq d \|\nabla f(x)\|^2$$

Randomized:

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} (\nabla_{i_t} f(x^t))^2$$

$$E_{i_t} [f(x^{t+1})] \leq E_{i_t} [f(x^t) - \frac{1}{2L} (\nabla_{i_t} f(x^t))^2]$$

$$\Rightarrow E_{i_t} [f(x^t)] - E_{i_t} \left[\frac{1}{2L} (\nabla_{i_t} f(x^t))^2 \right]$$

$$\Rightarrow f(x^t) - \frac{1}{2L} \sum_{i=1}^d p(i_t = i) (\nabla_{i_t} f(x^t))^2$$

$$\Rightarrow f(x^t) - \frac{1}{2L} \sum_{i=1}^d \left(\frac{1}{d}\right) (\nabla_{i_t} f(x^t))^2$$

$$\Rightarrow f(x^t) - \frac{1}{2Ld} \sum_{i=1}^d (\nabla_i f(x^t))^2 = f(x^t) - \frac{1}{2Ld} \|\nabla f(x)\|^2$$

$$p(i_t = i) = \frac{1}{d}$$

$$E_{i_t} [g(i_t)] = \sum_{i=1}^d p(i_t = i) g(i_t)$$

$$E(a+b) = E(a) + E(b)$$

$$E(c(a+b)) = cE(a) + cE(b)$$

"linearity of expectation"

$$E_{i+}[f(x^t)] \leq f(x^t) - \frac{1}{2L_i} \|\nabla f(x^t)\|^2$$

combine with lower bound as before, then we're done

Lipschitz $\nabla_i^2 f(x) \leq L_i$

$$p(i_+ = i) = \frac{L_i}{\sum_j L_j}$$

$$x^{t+1} = x^t - \frac{1}{L_i} \nabla f(x^t) e_{i+}$$

L_{i+} in upper bound

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{2L_{i+}} (\nabla_{i+} f(x^t))^2$$

$$1 - \frac{\mu}{2L} \Rightarrow 1 - \frac{\mu}{\sum L_i} \quad \begin{matrix} \sum L_i \leq d \max\{L_i\} \\ \sum L_i \leq dL \end{matrix}$$

Q4

final result:

$$\|x^{t+1} - x^*\| \leq \rho \|x^t - x^*\|$$

at x^t , $\rho < 1$, we are moving closer to solution

$$\|x^{t+1} - x^*\|^2 = \|x^t - x^*\|^2 - 2\alpha \nabla f(x^t)^\top (x^t - x^*) + \alpha^2 \|\nabla f(x^t)\|^2$$

$$- (\nabla f(x) - \nabla f(y))^\top (x - y) \leq \underbrace{-\frac{\mu L}{\mu + L} \|x - y\|^2}_{\substack{y=x^* \\ x=x^*}} - \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|^2$$

$(\nabla f(x^t) - \nabla f(x^*))$
 $\| \cdot \| \geq 0$ at optimal

$$= \|x^t - x^*\|^2 - 2\alpha \left[\frac{\mu L}{\mu + L} \|x^t - x^*\|^2 - \frac{1}{\mu + L} \underbrace{\|\nabla f(x^t) - \nabla f(x^*)\|^2}_{\text{zero}} \right] + \alpha^2 \|\nabla f(x^t)\|^2$$

$$= \left(1 - \frac{2\alpha\mu L}{\mu + L}\right) \|x^t - x^*\|^2 + \dots - \alpha \|\nabla f(x^t)\|^2 \left(\frac{2}{\mu + L} - \alpha\right)$$

goes to zero \downarrow
0 $\alpha = \frac{2}{\mu + L}$

4.2

step 1. $x^{t+1} = \text{prox} [x^t - \alpha_t \nabla f(x^t)]$

$$x^* = \text{prox} [x^* - \alpha_t \nabla f(x^*)]$$

step 2. $\|\text{prox}(x) - \text{prox}(y)\| \leq \|x - y\|$

continue as in 4.1, will give exact same as 4.1 in the end.

4.3

$$\log(\epsilon) = -\log(1/\epsilon)$$