

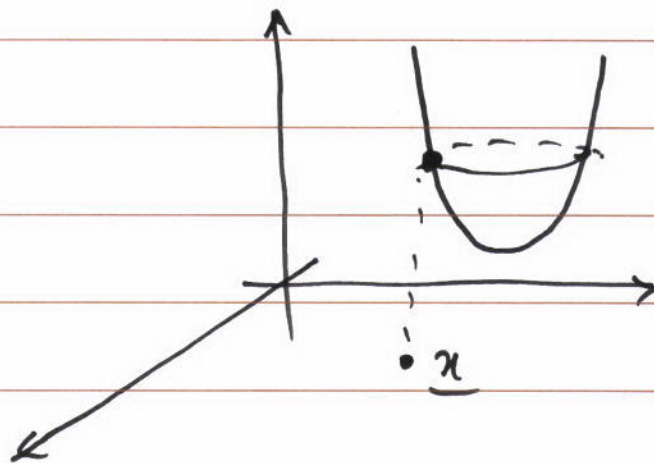
Discussion 4:

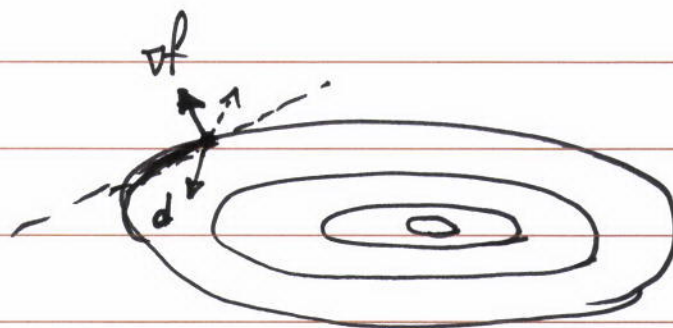
- Briefly talk about Iterative Descent Methods
 - preprocessing data
 - some tips for HW:
-

Q: How are we gonna achieve opt. point?

A: Iterative Descent Methods

- $\nabla f(x) = 0 \Rightarrow$ we have candidate
- $\nabla f(x) \neq 0 \Rightarrow$ NOT a candidate





if angle Δ (between ∇f and \underline{d} $> 90^\circ$)
 then direction \underline{d} results in a decrease
 in f

if $\nabla f^T \underline{d} < 0$

$\exists \delta > 0$ s.t. $f(\underline{x} + \alpha \underline{d}) < f(\underline{x}) \quad \forall \alpha \in (0, \delta)$

Iterative Descent Methods:

$$\underline{x}^{r+1} = \underline{x}^r + \underbrace{\alpha^r}_{\text{step size}} \underbrace{\underline{d}^r}_{\text{direction}}$$

Choices of direction \underline{d}^r :

- Steepest / Gradient Descent: $\underline{d}^r = -\nabla f(\underline{x}^r)$

- Diagonal scaled Gradient Descent:

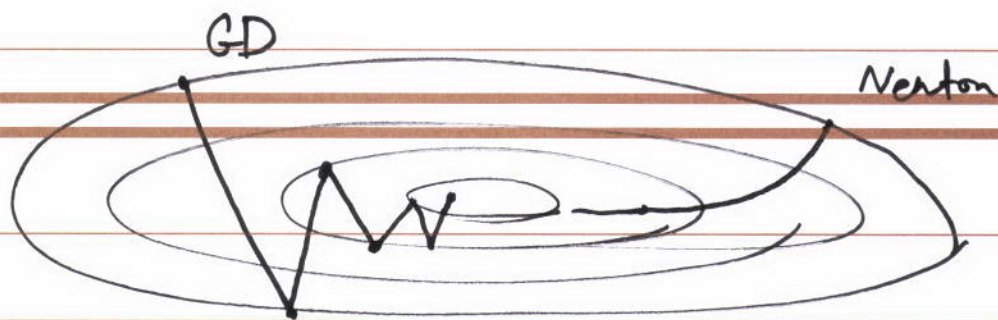
$$\underline{d}^r = -\underbrace{D^r}_{\text{diagonal}} \nabla f(\underline{x}^r)$$

$D^r \succeq 0$

• Newton Direction $\underline{d}^r = - (\nabla^2 f(\underline{x}^r))^{-1} \nabla f(\underline{x}^r)$

$$\underbrace{f(\underline{x}^r + \underline{d})}_{\min_{\underline{d}} f(\underline{x}^r + \underline{d})} \simeq f(\underline{x}^r) + \underbrace{\langle \underline{d}, \nabla f(\underline{x}^r) \rangle}_{\triangleq \nabla^T f(\underline{x}^r) \underline{d}} + \frac{1}{2} \underline{d}^T \nabla^2 f(\underline{x}^r) \underline{d}$$

$$\begin{aligned} \frac{\partial}{\partial \underline{d}} = 0 &\Rightarrow \nabla f(\underline{x}^r) + \nabla^2 f(\underline{x}^r) \underline{d}^* = 0 \\ &\Rightarrow \underline{d}^* = - (\nabla^2 f(\underline{x}^r))^{-1} \nabla f(\underline{x}^r) \end{aligned}$$



- Benefit : fast convergence
- Drawback : compute inverse

Diagonalized version of Newton:

$$\begin{aligned} D^r &= \text{diag} \left(\nabla^2 f(\underline{x}^r) \right)^{-1} \\ \underline{d}^r &= - D^r \nabla f(\underline{x}^r) \end{aligned}$$

Choice of step size:

• Constant $\alpha^r = \alpha \quad \forall r=0,1,\dots$ (prob: we don't exact const.)

• Exact Minimization $\alpha^r = \underset{\alpha \geq 0}{\operatorname{argmin}} f(\underline{x}^r + \alpha \underline{d}^r) \quad (1)$

• Limited Minimization $\alpha^r = \underset{\alpha \in (0, \bar{\alpha})}{\operatorname{argmin}} f(\underline{x}^r + \alpha \underline{d}^r)$

Adv.: — less computation compared to (1)

— less amount time compared to (1)

— we are NOT sure about opt. direction, so don't search entire space

• Diminishing $\alpha^r \downarrow 0$ with $\sum \alpha^r = \infty$

why? \downarrow need to search the entire space to reach to the opt.

$$\alpha^r = \frac{1}{r} \quad \checkmark \quad \sum \frac{1}{r} \rightarrow \infty$$

Convergence Rate:

— Asymptotic behavior of function f :

$$e(x) = \|x - x^*\| \quad \text{or} \quad e(x) = |f(x) - f(x^*)|$$

$$\lim_{r \rightarrow \infty} \sup \frac{e(x^{r+1})}{e(x^r)} = \beta$$

$$\frac{\|x^{r+1} - x^*\|}{\|x^r - x^*\|}$$

$$\beta = \begin{cases} \in (0, 1) & \text{linear} \rightarrow \frac{\|x^{r+1} - x^*\|}{\|x^r - x^*\|} = \frac{1}{2} \\ 1 & \text{sublinear} \\ 0 & \text{superlinear} \end{cases}$$

Example :

(for linear) $\Rightarrow \|x^{r+1} - x^*\| = \left(\frac{1}{2}\right)^r \|x^0 - x^*\|$

$\beta = \frac{1}{2}$

The smaller ^{the} β is, the faster convergence
* rate is.

- So far we saw GD alg.

- $\begin{cases} \text{Batch GD} \\ \text{Stochastic GD} \end{cases}$

Let's consider the following example to see these two GD (Batch & Stoch.):

Perceptron Problem:

- Given $\{(x_i, y_i)\}_{i=1}^N$ and having a $\{ \pm 1 \}$

linear model $g(w, x) = w^T x$ with hypothesis

$$h(x) = \text{sign}(w^T x)$$

- A prediction is correct if $y w^T x > 0$

Goal: try to find a "good" model s.t.

$h(x)$ makes few mis-prediction.

Formalize as optimization Prob. $\begin{cases} \text{- training data} \\ \text{- loss function} \end{cases}$

$$J(w) = \frac{1}{N} \sum_{i=1}^N L(w^T x_i, y_i)$$

Algorithm to find opt. w :

- initialize w
- Compute $\nabla J(w)$
- update $w^{(t+1)} = w^{(t)} - \eta \nabla J(w)$
step size
- repeat this process until convergence

\Rightarrow Batch GD (use all data to compute)

what loss function?

- 0/1 Loss function

$$J_{0/1}(w) = \frac{1}{N} \sum L(\text{sgn}(w^T x_i), y_i)$$

$$L(y_i', y_i) \triangleq \begin{cases} 0 & \text{if } y_i' = y_i \\ 1 & \text{else} \end{cases}$$

~~Not~~ (Not a good loss function because of flat shape of $J(w)$)

- perceptron criterion
(hinge loss)

$$J_p(w) = \frac{1}{N} \sum_{i=1}^N \max(0, -y_i w^T x_i)$$

\rightarrow is 0 if y_i is correctly predicted, otherwise is "confidence" in the misprediction.

In order to update \underline{w} , we don't need to have ALL data points, we can update \underline{w} by having every single data point.

for sample i : $J_i(\underline{w}) = \max(0, -y_i \underline{w}^T \underline{x}_i)$

$$\frac{\partial J_i}{\partial w_j} = \begin{cases} 0 & \text{if } y_i \underline{w}^T \underline{x}_i > 0 \\ -y_i x_{ij} & \text{else} \end{cases}$$

Vector \curvearrowright

$$\nabla J_i = \begin{cases} 0 & \text{if } y_i \underline{w}^T \underline{x}_i > 0 \\ -y_i \underline{x}_i & \text{else} \end{cases}$$

after observing (\underline{x}_i, y_i) , if it is mis predicted, we update \underline{w} as follows

$$\underline{w} \leftarrow \underline{w} + y_i \underline{x}_i$$

Algorithm for Stochastic GD:

- initialize \underline{w}

- Repeat

$$u_i \leftarrow \underline{w}^T \underline{x}_i$$

if $y_i u_i \leq 0$

$$\underline{w} \leftarrow \underline{w} + y_i \underline{x}_i$$

* η : (is also referred as learning rate)
(as learning continues, η should be smaller)

- Preprocessing of Data

- Represent Data in appropriate way

- Categorical
- Real
- ~~Real~~ Integer

- Standardizing data

- subtract by ~~mean~~ mean
- divide by standard deviation

- Handling missing data

Categorical Data

index	feature label	
1	A	1
2	B	1
3	C	1
4	F	0

(Encoding)

~~Converting~~ Converting Dummy variable V_i :

- for D distinct component, we are gonna ~~have~~ add D features which are 1 or 0

New form

index	f_1	f_2	f_3	f_4	label
1	1	0	0	0	1
2	0	1	0	0	1
3	0	0	1	0	1
4	0	0	0	1	0

4 distinct letters

A \rightarrow 1 0 0 0

B \rightarrow 0 1 0 0

C \rightarrow 0 0 1 0

F \rightarrow 0 0 0 1

drawback is if you have large D , you need to add so many features!!

• Standardizing Data:

(Variance) Assume feature 1 component over all samples are as follows

f_1
1
2
3

if we replace f_1 by f'_1 where

f'_1
4
8
12

should the model change weights for the 1st feature?

Since values in f'_1 are larger, model considers larger weights w'_1 compared to weight for f_1 , i.e. w_1 .

* To resolve this issue, we need to normalize all features by dividing all features by their standard deviation.

(Mean) Assume f_1 as before but f_1' as $\frac{f_1'}{1001}$
1001
1002
1003

if we replace f_1 by f_1' , should the model change the weights for the 1st feature?

Again, since values in f_1' is larger compared to f_1 , the weights for f_1' , i.e. w_1' , would be larger compared to the weight for f_1 .

* To resolve this issue, we need to subtract all features by their corresponding mean.

Therefore, by performing the aforementioned process all features have the same importance.

(Next Session for handling missing data)

Tips for HW:

python users (only use `numpy.shuffle` to shuffle the ~~into~~ indices, then code on your own to make cross-validation)