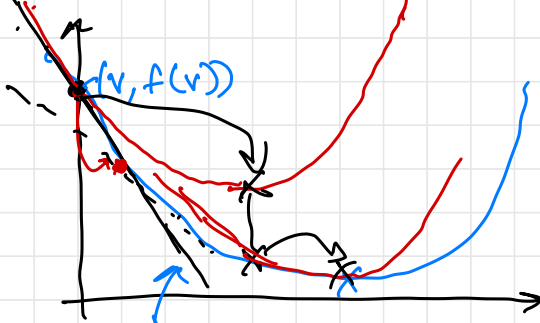


# 1/24/2023: Second-order optimization, Softmax Regression

Review: Gradient Descent = First-order optimization  
only use first derivative



First-order Taylor expansion  
 $y = f(v) + (x-v)f'(v)$

Q: Why must we take small steps?

A: First-order approx is not that reliable when taking large step

Q: Can we make a better approximation?

Second-order Taylor expansion

$$y = f(v) + (x-v)f'(v) + \frac{1}{2}(x-v)^2 f''(v)$$

Advantage: This has a global minimum if  $f''(v) > 0$

Algorithm (Newton-Raphson method in 1D)

Guess a point  $x^{(0)}$  (eg. 0)

for  $t=1, \dots, T$

Compute second-order approx. to  $f$  at  $x^{(t-1)}$

$$x^{(t)} \leftarrow \text{minimizer of this} \quad x^{(t)} \leftarrow x^{(t-1)} - \frac{f'(x^{(t-1)})}{f''(x^{(t-1)})}$$

return  $x^{(T)}$

take derivative, set equal to 0

$$f'(v) + (x-v) \cdot f''(v) = 0$$

$$x = v - \frac{f'(v)}{f''(v)}$$

# Newton-Raphson in $\mathbb{R}^d$

- Second derivative?
- Second-order Taylor expansion?
- Minimize  $\uparrow$

Hessian Matrix: For a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,

the Hessian  $H(f)$  is a  $d \times d$  matrix

$$\text{where } H(f)_{ij} = \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(x)$$

$d=2$ : 
$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial}{\partial x_1} \frac{\partial f}{\partial x_2} \\ \frac{\partial}{\partial x_1} \frac{\partial f}{\partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

Example

$$f(x) = x_1^3 + 5x_1^2 x_2$$

$$\frac{\partial f}{\partial x_1} = 3x_1^2 + 10x_1 x_2$$

$$\frac{\partial f}{\partial x_2} = 5x_1^2$$

$$\frac{\partial^2 f}{\partial x_1^2} = 6x_1 + 10x_2$$

$$\frac{\partial^2 f}{\partial x_2^2} = 0$$

$$\frac{\partial}{\partial x_1} \frac{\partial f}{\partial x_2} = 10x_1$$

$$H(f) = \begin{bmatrix} 6x_1 + 10x_2 & 10x_1 \\ 10x_1 & 0 \end{bmatrix}$$

# Second-order Taylor Expansion in $\mathbb{R}^d$

$$\text{Let } g = \nabla_x f(v), \quad H = H(f)(v)$$

$$f(x) \approx \underbrace{f(v) + g^T(x-v)}_{\text{first-order approx}} + \frac{1}{2} (x-v)^T H (x-v)$$

In general, the expression  $u^T A u$  is called a quadratic form (general form of a quadratic)

$$u^T A u = \sum_{i=1}^d \sum_{j=1}^d A_{ij} u_i u_j \quad \leftarrow \nabla_u u^T A u = 2 A u$$

minimizing the Taylor expansion

→ compute gradient, set equal to 0

$$g + H(x-v) = 0$$

$$x = v - H^{-1} g$$

Update rule for Newton-Raphson in  $\mathbb{R}^d$

## Newton-Raphson vs. Gradient Descent

↓  
 $O(d^2)$  memory to store H

↓  
 $O(d)$  memory

Second-order methods expensive when  $d$  is large

In practice: use a low-rank approximation to H that takes  $O(d)$  memory

L-BFGS: Approximate H + conservative update

↑ Limited Memory

— Name of algorithm

# Announcements

HW0 — Solutions

Grades by thurs

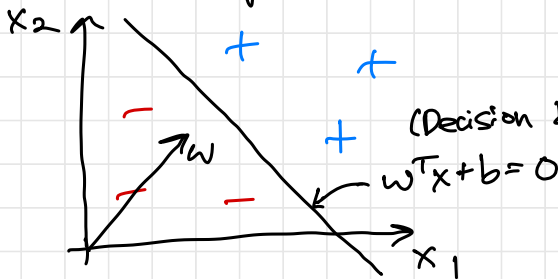
Section Fri

HW1 — out tonight (due Feb 7)

## Softmax Regression ("Multinomial Logistic Regression")

Multi-class classification:  $C$  different classes

- Handwritten digit  $\rightarrow [0, 1, 2, \dots, 9]$  ( $C=10$ )
- Image  $\rightarrow [bird, mammal, reptile, \dots]$

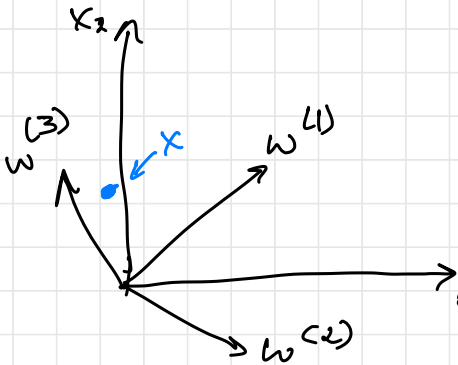


predict 1 if  $w^T x + b > 0$   
-1 else

dot product with  $w$   
measures how much you  
look like class 1

Softmax Regression: parameter vectors  
 $w^{(1)}, \dots, w^{(c)} \in \mathbb{R}^d$   
 total  $C \times d$  parameters

$w^{(j)T} x$  measures how much  $x$  looks like class  $j$



Decision Rule: given  $x$ ,  
 compute  $w^{(1)T} x, \dots, w^{(c)T} x$   
 return  $j$  with largest  $w^{(j)T} x$

Probabilistic Story:

$$p(y=j | x; w) = \frac{\exp(w^{(j)T} x)}{\sum_{k=1}^c \exp(w^{(k)T} x)}$$

$w^{(1)T} x = 1 \rightarrow \exp \approx 2.7 \rightarrow p(y=1 | x; w) \approx 0.27$   
 $w^{(2)T} x = -3 \rightarrow \exp \approx 0.1 \rightarrow p(y=2 | x; w) \approx 0.01$   
 $w^{(3)T} x = 2 \rightarrow \exp \approx 7.4 \rightarrow p(y=3 | x; w) \approx 0.72$   
 Sum  $\approx 10.2$

Maximum Likelihood Estimation

Softmax function

$$NLL(w) = - \sum_{i=1}^n \log P(y=y^{(i)} | x^{(i)}; w)$$

"negative log likelihood"  
 want to minimize

$$= - \sum_{i=1}^n w^{(y^{(i)})T} x^{(i)} - \log \left( \sum_{k=1}^c \exp(w^{(k)T} x^{(i)}) \right)$$

our loss function!

Gradient time:

$$\nabla_{w^{(j)}} NLL(w) = - \sum_{i=1}^n \mathbb{I}[y^{(i)}=j] \cdot x^{(i)}$$

(Do this for every  $j$ )

$$- \frac{1}{\sum_{k=1}^c \exp(w^{(k)T} x^{(i)})} \cdot \exp(w^{(j)T} x^{(i)}) \cdot x^{(i)}$$

scalar  $p(y=j | x^{(i)}; w)$

$$= \sum_{i=1}^S \left( \underbrace{p(y=j | x^{(i)}; \omega)}_{\text{between 0 \& 1}} - \mathbb{1}[y^{(i)}=j] \right) x^{(i)}$$

IF  $y^{(i)} \neq j$  then positive  $\cdot x^{(i)} \Rightarrow$  GD subtracts multiple of  $x^{(i)}$

$y^{(i)} = j$  then negative  $\cdot x^{(i)} \Rightarrow$  GD add multiple of  $x^{(i)}$

$$w^{(t)} \leftarrow w^{(t-1)} - \underbrace{\gamma \cdot \nabla f(w^{(t-1)})}$$