

Lecture 7: 2/18/17

Unconstrained Optimization

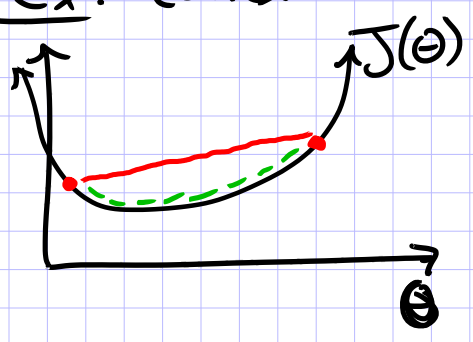
Given function $J(\theta)$, $J: \mathbb{R}^M \rightarrow \mathbb{R}$

Goal is $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$ by convention
 $= \underset{\theta}{\operatorname{argmax}} -J(\theta)$

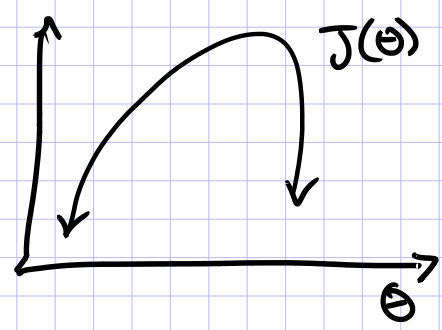
For ML: J is usually an objective fn.
 θ is vector of parameters

Convexity

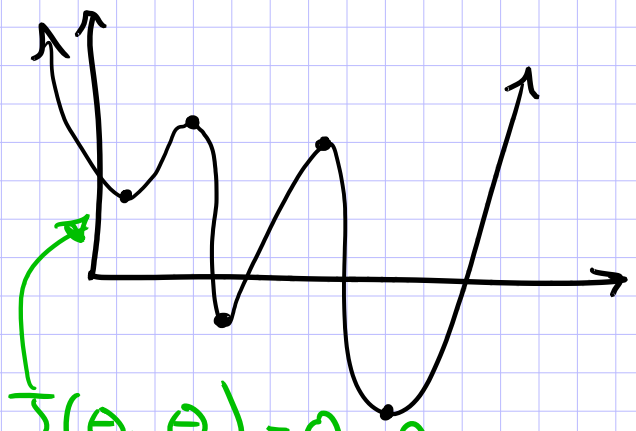
Ex: convex



Ex: concave



Ex: nonconvex



Ex: $J(\theta_1, \theta_2) = \theta_1 \cdot \theta_2$

Def: The gradient of J

$$\nabla J(\theta) \triangleq \begin{bmatrix} \frac{dJ(\theta)}{d\theta_1} \\ \frac{dJ(\theta)}{d\theta_2} \\ \vdots \\ \frac{dJ(\theta)}{d\theta_M} \end{bmatrix}$$

first order partial derivatives

Def: The Hessian matrix: second-order partial derivs

$$\nabla^2 J(\vec{\theta}) \triangleq H_J(\vec{\theta}) \triangleq \begin{bmatrix} \frac{d^2 J(\vec{\theta})}{d\theta_1^2} & \dots & \frac{d^2 J(\vec{\theta})}{d\theta_M d\theta_1} \\ \vdots & & \vdots \\ \frac{d^2 J(\vec{\theta})}{d\theta_1 d\theta_M} & \dots & \frac{d^2 J(\vec{\theta})}{d\theta_M^2} \end{bmatrix}$$

$H_{ij} = \frac{d^2 J(\vec{\theta})}{d\theta_i d\theta_j}$

Def: $J'(\theta) \triangleq \frac{dJ(\theta)}{d\theta}$
 $J''(\theta) \triangleq \frac{d^2 J(\theta)}{d\theta^2}$

Closed-Form Opt.

Ex #1: 1D function $J(\theta)$

① Solve $J'(\theta) = 0$ for a point $\hat{\theta}$

② Q: Is $\hat{\theta}$ a min, max, or saddle point?

A: If $J''(\theta) > 0 \Rightarrow \text{min}$

If $J''(\theta) < 0 \Rightarrow \text{max}$

If $J''(\theta) = 0 \Rightarrow$ use "extremum test" (higher-order part)

← second derivative test.

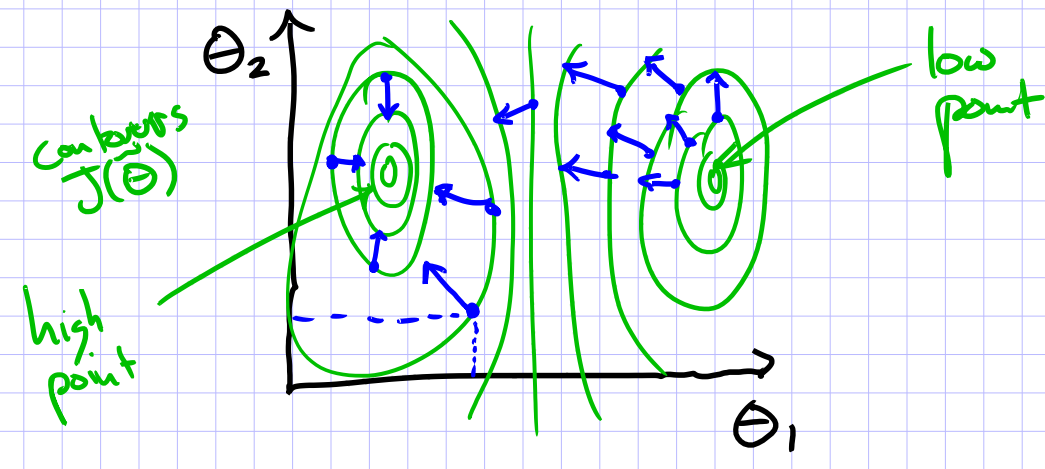
Ex #2: M-dim function $J(\vec{\theta})$, $\vec{\theta} \in \mathbb{R}^M$

test for min, max, s.p. requires Hessian

Gradient Descent

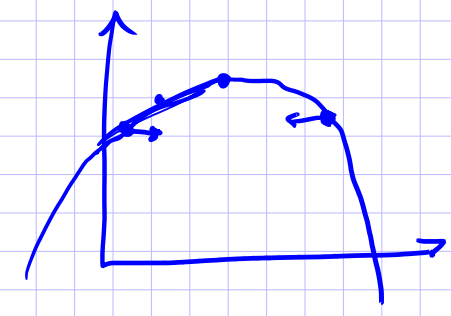
- Motiv: - What if we can't solve $\nabla J(\vec{\theta}) = \vec{0}$ analytically? (e.g. Logistic Regression)
- What if we can, but it's inefficient to compute. (e.g. Linear Reg.)

Ex: Gradients in 2D



$\nabla J(\vec{\theta}) = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$

blue arrows indicate the direction of the gradient, but gradients also have magnitude



$\nabla J(\theta) = [\]$

Algorithm:

- 1) Choose an initial **starting point** $\vec{\theta}$
- 2) Repeat
 - a) Compute the gradient $\nabla J(\vec{\theta})$
 - b) Choose a **step size** size γ
 - c) Update $\vec{\theta} = \vec{\theta} - \gamma \nabla J(\vec{\theta})$
- 3) Return $\vec{\theta}$ when **stopping criterion** satisfied.

Starting Point:

- Two good options for ML:
- (1) $\vec{\theta} = 0$
 - (2) $\vec{\theta}$ randomly

Stopping Criterion

Typically

$$\| \nabla J(\vec{\theta}) \|_2 < \epsilon$$

$$\| \vec{u} \|_2 \triangleq \sqrt{\sum_{i=1}^n u_i^2}$$

Step Size:

- Lots of Options
- ~~fixed value eg. $\gamma = 0.1$~~
 - Exact line search
 - Backtracking line search

Stochastic Gradient Descent

Assume: $J(\vec{\theta}) \triangleq \frac{1}{N} \sum_{i=1}^N J_i(\vec{\theta})$ where $J_i: \mathbb{R}^M \rightarrow \mathbb{R} \forall i$

Ex: $p(\mathcal{D}|\theta) = \prod_{i=1}^N p(x^{(i)}, y^{(i)}|\theta)$

$$l(\theta) = \log p(\mathcal{D}|\theta) = \sum_{i=1}^N \underbrace{\log p(x^{(i)}, y^{(i)}|\theta)}_{l_i(\theta)}$$

neg. log. like $\rightarrow J(\theta) = -\frac{1}{N} l(\theta) = \frac{1}{N} \sum_{i=1}^N \underbrace{-l_i(\theta)}_{J_i(\theta)}$

Expectations of Gradients

$$\frac{dJ(\vec{\theta})}{d\theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\theta_j} (J_i(\vec{\theta}))$$

$$\nabla J(\vec{\theta}) = \begin{bmatrix} \vdots \\ \text{jth} \\ \vdots \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$

Recall: for any discrete r.v. X $E_X[f(X)] \triangleq \sum_x P(X=x) f(x)$

Q: What is the expected value of a randomly chosen $\nabla J_i(\theta)$?

Let $I \sim \text{Uniform}(\{1, \dots, N\})$

$$\Rightarrow P(I=i) = \frac{1}{N} \text{ if } i \in \{1, \dots, N\}$$

$$E_I[\nabla J_I(\vec{\theta})] = \sum_{i=1}^N P(I=i) \nabla J_i(\vec{\theta})$$

$$= \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$

$$= \nabla J(\vec{\theta})$$

SGD Algo:

- 1) Choose **starting point** $\vec{\theta}$
- 2) For t in $1 \dots T$
 - a) sample $i \sim \text{Unif}(\{1, \dots, N\})$
 - b) compute $\vec{d} = \nabla J_i(\vec{\theta})$ ← stochastic approx. to true gradient $\nabla J(\vec{\theta})$
 - c) choose **step size** $\gamma^{(t)}$
 - d) update $\vec{\theta} = \vec{\theta} - \gamma^{(t)} \vec{d}$
- 3) Return $\vec{\theta}$ when **stopping criteria** satisfied

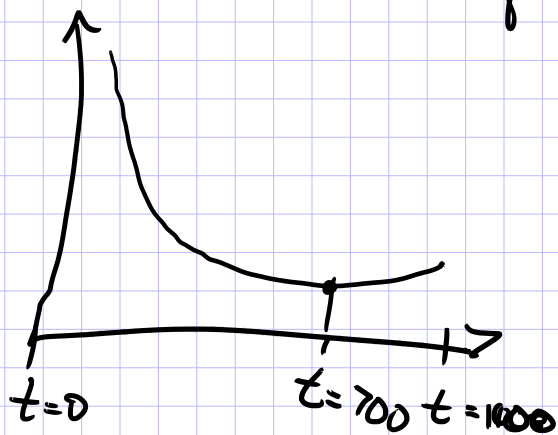
Step Size:

- larger values early on smaller values closer to opt.
- heuristic:
 - ① pick an initial $\gamma = 0.1$
 - ② run until no improvement
 - ③ $\gamma = \gamma/2$
- learning rate schedule

$$\gamma^{(t)} = \frac{\gamma^{(0)}}{(t-1)\gamma^{(0)} + 1}$$

Stopping Criteria:

For ML:
plot validation error and choose lowest point



true unknown min

Convergence Analysis:

Def: convergence is when $J(\vec{\theta}) - J(\vec{\theta}^*) < \epsilon$

Methods	Steps to Converge	Computation per iteration
Newton's Method	$O(\ln \ln 1/\epsilon)$	$\nabla^2 J(\theta)$ $\nabla^2 J(\theta) \leftarrow O(NM^2)$
GD	$O(\ln 1/\epsilon)$	$\nabla J(\theta) \leftarrow O(NM)$
SGD	$O(1/\epsilon)$	$\nabla J_i(\theta) \leftarrow O(M)$

"almost sure" convergence lots of caveats and conditions

way less computation

Takeaway: SGD has much slower asymptotic convergence. but is often faster in practice.

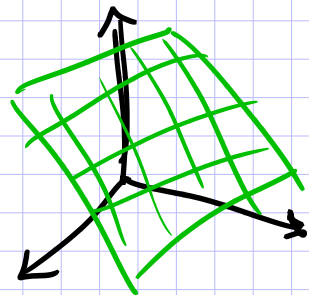
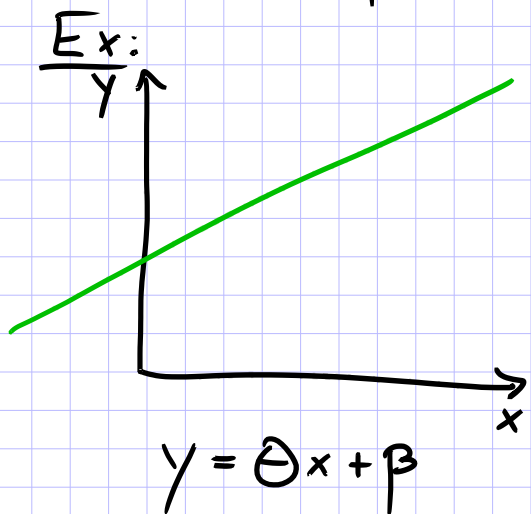
Regression Problem

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$$

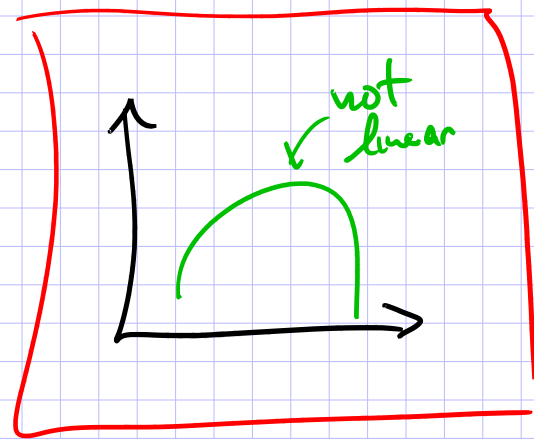
$$\vec{x} \in \mathbb{R}^M \quad \leftarrow \text{input, features}$$

$$y \in \mathbb{R} \quad \leftarrow \text{output}$$

Linear Function

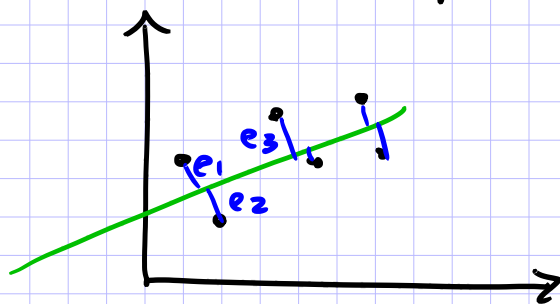


$$y = \vec{\theta}^T \vec{x} + \beta$$
$$\vec{\theta} \in \mathbb{R}^M, \beta \in \mathbb{R}$$



Def: residuals: distance from observed value to predicted value

$$e_i \triangleq y^{(i)} - (\vec{\theta}^T \vec{x} + \beta)$$



Notation Trick:

Define $x_0 \triangleq 1$

$$y = \vec{\theta}^T \vec{x} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_M \end{bmatrix}^T \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_M \end{bmatrix}$$

replacement for β

always 1

Linear Regression as Functn Approx.

Story:

① Assume D generated: $x^{(i)} \sim p^*(x)$
 $y^{(i)} = h^*(x^{(i)})$ unkn.

② Choose hypothesis space, H

$$H = \{h(\vec{x}) = \vec{\theta}^T \vec{x} : \vec{\theta} \in \mathbb{R}^M\}$$

↑ space of all linear fun. in M -dim.

③ Choose an objective fun.

Minimize mean squared error (MSE)

$$J(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N e_i^2 = \frac{1}{N} \sum_{i=1}^N (h^*(x^{(i)}) - h(x^{(i)}))^2$$

④ Solve unconstrained opt. problem. via favorite method:

$$\hat{\vec{\theta}} = \underset{\vec{\theta}}{\operatorname{argmin}} J(\vec{\theta})$$

- closed-form ← "Normal Eqs"
- SGD ← "Least Mean Squares"
- GD
- Newton

⑤ Predict: $\hat{y} = h(\vec{x}) = \hat{\vec{\theta}}^T \vec{x}$