

Recitation 2

Gradient Descent and Stochastic Gradient Descent

DS-GA 1003 Machine Learning

Spring 2021

Feburary 10, 2021

Agenda

- Gradient Descent
 - Adaptive Learning Rate
 - Coding Exercise
- Stochastic Gradient Descent
 - Coding Exercise
- Application
 - Linear Regression
 - Logistic Regression

Gradient Descent

Gradient Descent Recap

Gradient Descent

- Initialize $x = 0$

- Repeat:

$$\bullet x \leftarrow \underline{x} - \eta \underline{\nabla f(x)}$$

step size

$$\begin{array}{c} \text{parameters} \\ \downarrow \\ \nabla f(x) \end{array}$$

direction
 $\|\nabla f(x)\| -$
 $\nabla f(x)$, gradient operator

- Until stopping criterion satisfied

- Choosing the step size is the key in gradient descent
- A fixed step size will work, eventually, as long as it's small enough

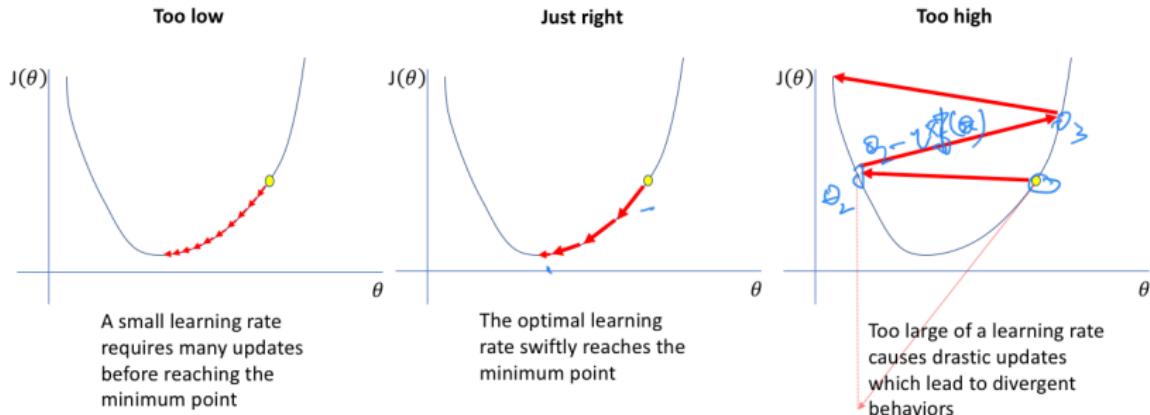
Gradient Descent

Gradient Descent Algorithm

- Goal: find $\theta^* = \arg \min_{\theta} J(\theta)$ *loss for*
- $\theta^0 := [\text{initial condition}]$
- $i := 0$
- while not [termination condition]:
 - compute $\nabla J(\theta_i)$ *gradient*
 - $\epsilon_i := [\text{choose learning rate at iteration } i]$
 - $\theta^{i+1} := \theta^i - \epsilon_i \nabla J(\theta_i)$
 - $i = i + 1$
- return θ^i

Gradient Descent

- How to initialize θ^0 ?
 - sample from some distribution
 - compose θ^0 using some heuristics
- + Normal / Gaussian
+ Uniform
+ Initialize to Zeros
- How to choose termination conditions?
 - run for a fixed number of iteration
 - the value of $f(\theta)$ stabilizes
 - θ^i converges
- Loss, Accuracy — Validation
- What is a good learning rate?



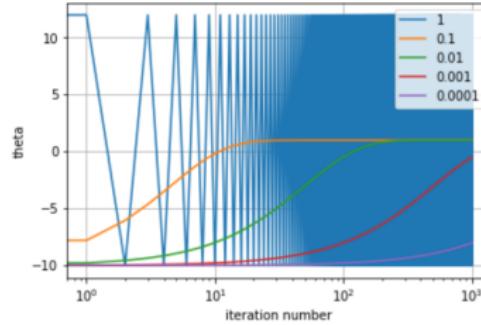
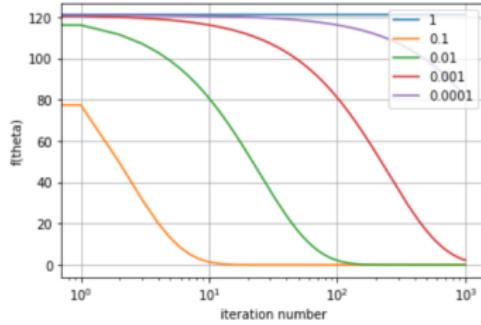
Learning Rate

Application

Suppose we would like to find $\theta^* \in \mathbb{R}$ that minimizes $f(\theta) = \theta^2 - 2\theta + 1$. The gradient (in this case, the derivative) $\nabla f(\theta) = 2\theta - 2$. We can easily see that $\theta^* = \operatorname{argmin}_\theta f(\theta) = 1$.

Learning Rate

- We applied gradient descent for 1000 iterations on $f(\theta) = \theta^2 - 2\theta + 1$ with varying learning rate $\epsilon \in \{1, 0.1, 0.01, 0.001, 0.0001\}$
- When the learning rate is too large ($\epsilon = 1$), $f(\theta)$ does not decrease through iterations. The value of θ_i at each iteration significantly fluctuates.
- When the learning rate is too small ($\epsilon = 0.0001$), $f(\theta)$ decreases very slowly.

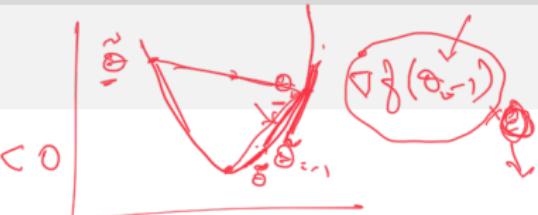


Adaptive Learning Rate

- Instead of using a fixed learning rate through all iterations, we can adjust our learning rate in each iteration using a simple algorithm.
- At each iteration i :
 - $\tilde{\theta} := \theta_{i-1} - \epsilon_{i-1} \nabla f(\theta_{i-1})$
 - $\delta := f(\theta_{i-1}) - f(\tilde{\theta})$
 - if $\delta \geq \text{threshold}$:
 - we achieve a satisfactory reduction on $f(\theta)$
 - $\theta_i = \tilde{\theta}$
 - maybe we can consider increasing the learning rate for next iteration
 $\epsilon_j := 2\epsilon_{i-1}$
 - else:
 - the reduction is unsatisfactory
 - $\theta_i = \theta_{i-1}$
 - the learning rate is too large, so we reduce the learning rate
 $\epsilon_i := \frac{1}{2}\epsilon_{i-1}$

Adaptive Learning Rate

$$f(\theta_{i-1}) - f(\tilde{\theta}) < 0$$



How to decide a proper threshold for $f(\theta_{i-1}) - f(\tilde{\theta})$

Armijo rule

If learning rate ϵ satisfies

$$f(\theta_{i-1}) - f(\tilde{\theta}) \geq \frac{1}{2} \epsilon \|\nabla f(\theta_{i-1})\|^2$$

s

how much change
gradient

then $f(\theta)$ is guaranteed to converge to a (local) minimum under certain technical assumptions.

You can find more details at [this link](#)

- 10 steps $\downarrow \epsilon$ by $1/2 \cdot 1/10$
- accuracy doesn't change
for x iterations

Stochastic Gradient Descent

Stochastic Gradient Descent

Stochastic Gradient Descent

- Initialize $w = 0$
- Repeat:
 - randomly choose training point $(x_i, y_i) \in \mathcal{D}_n$
 - $w \leftarrow w - \eta \underbrace{\nabla_w \ell(f_w(x_i), y_i)}_{\text{Grad(Loss on i'th example)}}$
- Until stopping criterion satisfied

- Equivalent to Minibatch Gradient Descent with batch size $N = 1$.
- Use a single randomly chosen point to determine step direction.

Minibatch Gradient Descent

Minibatch Gradient Descent (minibatch size N)

- Initialize $w = 0$
 - Repeat:
 - randomly choose N points $\{(x_i, y_i)\}_{i=1}^N \subset \mathcal{D}_n$
 - $w \leftarrow w - \eta \left[\frac{1}{N} \sum_{i=1}^N \nabla_w \ell(f_w(x_i), y_i) \right]$
 - Until stopping criterion satisfied
-
- Minibatch gradient is an unbiased estimate of full-batch gradient:
$$\mathbb{E} \left[\nabla \hat{R}_N(w) \right] = \nabla \hat{R}_n(w)$$
 - Use a random subset of size N to determine step direction
 - Bigger N = Better estimate of the gradient, but slower (more data to touch)
 - Smaller N = Worse estimate of the gradient, but faster

Application

Gradient Descent for Linear Regression

Linear Least Squares Regression Setup

- **Data:** Inputs are feature vectors of dimension d. Outputs are continuous scalars.

$$\mathcal{D} = \left\{ \mathbf{x}^{(i)}, y^{(i)} \right\}_{i=1}^n \text{ where } \mathbf{x} \in \mathbb{R}^d \text{ and } y \in \mathbb{R}$$

Input Space Action Space Output Space

- **Hypothesis Space:** $\mathcal{F} = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(\mathbf{x}) = \underline{\theta^T \mathbf{x}}, \theta \in \mathbb{R}^d\}$
- **Action:** Our prediction is a linear function of the inputs

$$\hat{y} = f_{\theta}(\mathbf{x}) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d \quad ?$$

$\hat{y} = f_{\theta}(\mathbf{x}) = \underline{\theta^T \mathbf{x}}$

(We assume x_1 is 1)

- **Loss:** $\ell(\hat{y}, y) = (y - \hat{y})^2$

Gradient Descent for Linear Regression

- **Goal:** Finding the set of parameters that minimize the empirical risk:

$$\hat{R}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\theta^T x^{(i)} - y^{(i)} \right)^2$$

where $\theta \in R^d$ parameterizes the hypothesis space \mathcal{F}

- Set our cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left(\theta^T x^{(i)} - y^{(i)} \right)^2$$

Three Approach to solving $\theta^* = \operatorname{argmin} J(\theta)$



- **Approach 1:** Closed Form Solution (set derivatives equal to zero and solve for parameters)
 - pros: one shot algorithm!
 - cons: does not scale to large datasets (matrix inverse is bottleneck)
- **Approach 2:** Gradient Descent (take larger, more certain steps toward the negative gradient)
 - pros: conceptually simple, guaranteed convergence
 - cons: batch, often slow to converge
- **Approach 3:** Stochastic Gradient Descent (take many small, quick steps opposite the gradient)
 - pros: memory efficient, fast convergence, less prone to local optima
 - cons: convergence in practice requires tuning and fancier variants

Approach 1: Close-Form Solution

- Transform the cost function in matrix form

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^n \left(x_i^T \theta - y_i \right)^2 = \frac{1}{n} \sum_{i=1}^n (\theta^T x_i - y_i)^2 \\ &= \frac{1}{2} (X\theta - \bar{y})^T (X\theta - \bar{y}) = \frac{1}{2} \|X\theta - y\|_2^2 \end{aligned}$$

- To minimize $J(\theta)$, take derivative and set to zero:

$$\nabla J(\theta) = (X^T X \theta - X^T y) = X^T (X\theta - y) = 0$$

$$\Rightarrow \hat{\theta} = (X^T X)^{-1} X^T y$$

- Ensure invertibility of $X^T X$
- What if X has less than full column rank?

Approach 2: Iterative Method GD

Gradient Descent Algorithm

- $\theta^0 := [\text{initial condition}]$
- $i := 0$
- while not [termination condition]:
 - $\theta^{i+1} := \theta^i - \epsilon_i \nabla_{\theta} J(\theta)$
 - $i = i + 1$
- return θ^i

Recall:

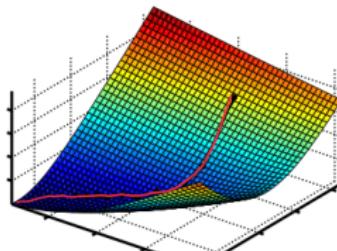
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_d} J(\theta) \end{bmatrix}$$

Approach 3: Iterative Method SGD

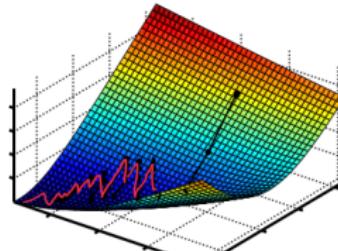
Stochastic Gradient Descent Algorithm

- $\theta^0 := [\text{initial condition}]$
- $i := 0$
- while not [termination condition]:
 - For each training pair $(x^{(j)}, y^{(j)})$ (in random order)
 - $\theta^{i+1} := \theta^i - \epsilon_i \nabla_{\theta} J^{(j)}(\theta)$ {with $J^{(j)}(\theta) = \frac{1}{2} (\theta^T x^{(j)} - y^{(j)})^2$ }
 - $i = i + 1$
- return θ^i

GD



SGD



Gradient Descent for Logistic Regression

Binary Classification Setup for Logistic Regression

- **Data:** Inputs are feature vectors of dimension d . Targets are class labels. $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^n$ where $x \in \mathbb{R}^d$ and $y \in \{0, 1\}$ — Outcome space
- **Action:** Our prediction is the probability of class label given linear signals



$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}} \quad \text{with} \quad g(z) = \frac{1}{1+e^{-z}}$$

- Sigmoid Function $g(z)$ takes a real-valued number and maps it into the range $[0, 1]$ (Probability Interpretation)
- Assume Action space

$$\begin{cases} P(y = 1 | x; \theta) = h_{\theta}(x) \\ P(y = 0 | x; \theta) = 1 - h_{\theta}(x) \end{cases}$$

- More Compactly: $p(y | x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$

Learning Logistic Regression

- Assumption: $(x_1, y_1), \dots, (x_n, y_n)$ are independently generated
- Likelihood:** The probability of getting the y_1, \dots, y_n in \mathcal{D} from the corresponding x_1, \dots, x_n

$$\begin{aligned} P(y_1, \dots, y_n | x_1, \dots, x_n) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n (h_\theta(x^{(i)}))^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

$\log(a \cdot b) = \log a + \log b$
 $\log x^n = n \log x$

- Goal:** maximize the log likelihood (Easier)

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$$

- Equivalent to minimize the objective function with logistic loss:

$$J(\theta) = \sum_{i=1}^n \ell(h_\theta(x^{(i)}), y^{(i)})$$

where $\ell(h(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$ - NLL



Learning Logistic Regression

- Analytic solution won't work
- Find optimum using iterative methods: Gradient Ascent or Stochastic Gradient Ascent
 - Gradient ascent rule: $\theta := \theta + \alpha \nabla_{\theta} \ell(\theta)$
 - Stochastic gradient ascent rule: $\theta := \theta + \alpha \nabla_{\theta} \ell^{(i)}(\theta)$ for random training pair (x^i, y^i)
- For one training example (x, y) , the partial derivative of log likelihood $\ell(\theta)$:

$$\frac{\partial}{\partial \theta} \ell(g(\theta)) = \frac{\partial \ell(g(\theta))}{\partial g(\theta)} \cdot \frac{\partial g(\theta)}{\partial \theta}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left(y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) g(\theta^T x)(1-g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1-g(\theta^T x)) - (1-y)g(\theta^T x))x_j \\ &= (y - h_{\theta}(x)) \cdot x_j \end{aligned}$$

References

- DS-GA 1003 Machine Learning Spring 2020
- **Stanford CS229 Note1**
- **CMU 10-701 Linear Regression Slide**