

Optimization for ML (3)

CS771: Introduction to Machine Learning

Piyush Rai

Stochastic Gradient Descent (SGD)

2

Writing as an average instead of sum.
Won't affect minimization of $L(\mathbf{w})$

- Consider a loss function of the form $L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ell_n(\mathbf{w})$

- The (sub)gradient in this case can be written as

$$\mathbf{g} = \nabla_{\mathbf{w}} L(\mathbf{w}) = \nabla_{\mathbf{w}} \left[\frac{1}{N} \sum_{n=1}^N \ell_n(\mathbf{w}) \right] = \frac{1}{N} \sum_{n=1}^N \mathbf{g}_n$$

Expensive to compute – requires
doing it for all the training
examples in each iteration ☹

(Sub)gradient of the loss
on n^{th} training example

- Stochastic Gradient Descent (SGD) approximates \mathbf{g} using a single training example
- At iter. t , pick an index $i \in \{1, 2, \dots, N\}$ uniformly randomly and approximate \mathbf{g} as

$$\mathbf{g} \approx \mathbf{g}_i = \nabla_{\mathbf{w}} \ell_i(\mathbf{w})$$

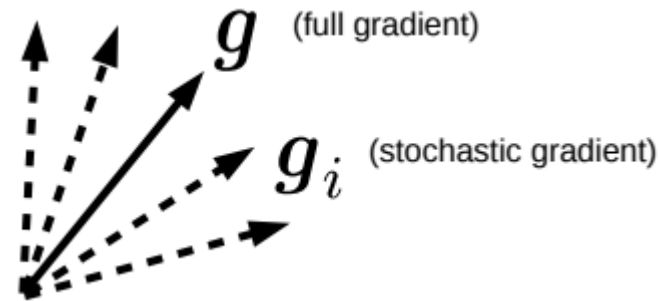
Can show that \mathbf{g}_i is an
unbiased estimate of \mathbf{g} ,
i.e., $\mathbb{E}[\mathbf{g}_i] = \mathbf{g}$

- May take more iterations than GD to converge but each iteration is much faster 😊
 - SGD per iter cost is $O(D)$ whereas GD per iter cost is $O(ND)$



Minibatch SGD

- Gradient approximation using a single training example may be noisy



The approximation may have a **high variance** – may slow down convergence, updates may be unstable, and may even give sub-optimal solutions (e.g., local minima where GD might have given global minima)

- We can use $B > 1$ unif. rand. chosen train. ex. with indices $\{i_1, i_2, \dots, i_B\} \in \{1, 2, \dots, N\}$
- Using this “minibatch” of examples, we can compute a minibatch gradient

$$\mathbf{g} \approx \frac{1}{B} \sum_{b=1}^B \mathbf{g}_{i_b}$$

- Averaging helps in reducing the variance in the stochastic gradient
- Time complexity is $O(BD)$ per iteration in this case



Constrained Optimization



Projected Gradient Descent

- Consider an optimization problem of the form

$$\mathbf{w}_{opt} = \arg \min_{\mathbf{w} \in \mathcal{C}} L(\mathbf{w})$$

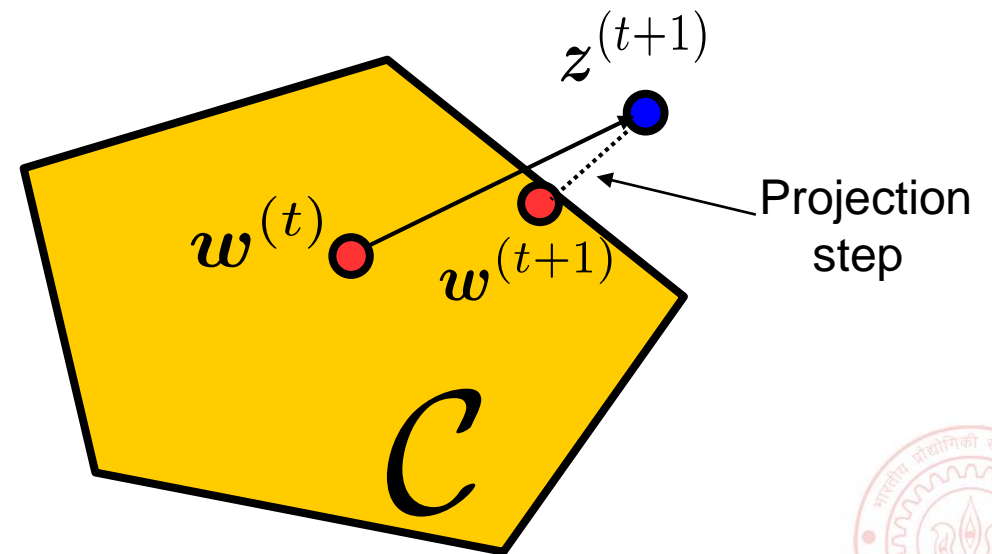
- Projected GD is very similar to GD with an extra **projection step**
- Each iteration t will be of the form

- Perform update: $\mathbf{z}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$

- Check if $\mathbf{z}^{(t+1)}$ satisfies constraints

- If $\mathbf{z}^{(t+1)} \in \mathcal{C}$, set $\mathbf{w}^{(t+1)} = \mathbf{z}^{(t+1)}$
- If $\mathbf{z}^{(t+1)} \notin \mathcal{C}$, project as $\mathbf{w}^{(t+1)} = \Pi_{\mathcal{C}}[\mathbf{z}^{(t+1)}]$

Projection
operator



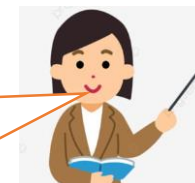
Projected GD: How to Project?

- Here projecting a point means finding the “closest” point from the constraint set

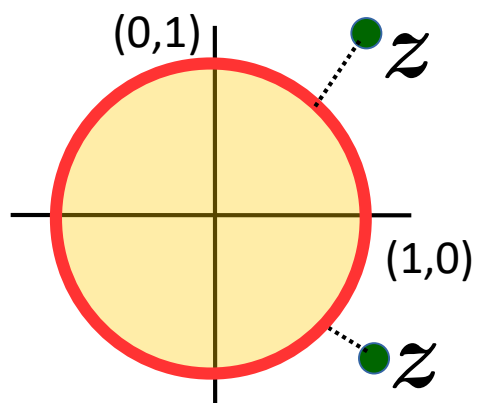
$$\Pi_{\mathcal{C}}[\mathbf{z}] = \arg \min_{\mathbf{w} \in \mathcal{C}} \|\mathbf{z} - \mathbf{w}\|^2$$

- For some sets \mathcal{C} , the projection step is easy

Projected GD commonly used only when the projection step is simple and efficient to compute



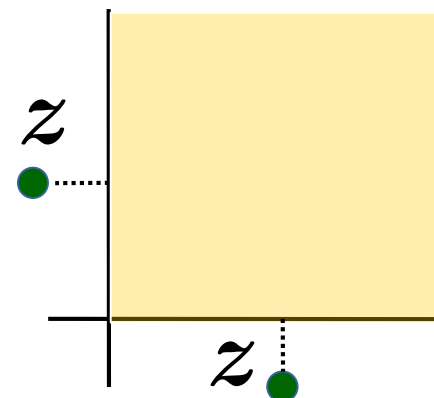
\mathcal{C} : Unit radius ℓ_2 ball



Projection = Normalize to unit Euclidean length vector

$$\hat{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_2 \leq 1 \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \text{if } \|\mathbf{x}\|_2 > 1 \end{cases}$$

\mathcal{C} : Set of non-negative reals



Projection = Set each negative entry in \mathbf{z} to be zero

$$\hat{\mathbf{x}}_i = \begin{cases} \mathbf{x}_i & \text{if } \mathbf{x}_i \geq 0 \\ 0 & \text{if } \mathbf{x}_i < 0 \end{cases}$$



Proximal Gradient Descent

- Consider minimizing a regularized loss function of the form

$$\arg \min_{\mathbf{w}} L(\mathbf{w}) + R(\mathbf{w})$$

Note: The reg. hyperparam. λ assumed part of $R(\mathbf{w})$ itself

- Proximal GD popular when regularizer $R(\mathbf{w})$ is non-differentiable
 - Basic idea: Do GD on $L(\mathbf{w})$ and use a prox. operator to regularize via $R(\mathbf{w})$
 - For a func. R , its prox. operator is $\text{prox}_R(\mathbf{z}) = \arg \min_{\mathbf{w}} \left[R(\mathbf{w}) + \frac{1}{2} \|\mathbf{z} - \mathbf{w}\|_2^2 \right]$
- Proximal GD**

- Assume reg. loss function of the form $L(\mathbf{w}) + R(\mathbf{w})$
- Initialize \mathbf{w} as $\mathbf{w}^{(0)}$
- For iteration $t = 0, 1, 2, \dots$ (or until convergence)

- Calculate the (sub)gradient of train. Loss (w/o reg.)

$$\mathbf{g}^{(t)} \in \partial L(\mathbf{w}^{(t)})$$

- Set learning rate η_t
- Step 1: $\mathbf{z}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$
- Step 2: $\mathbf{w}^{(t+1)} = \text{prox}_R(\mathbf{z}^{(t+1)})$

Special Cases

For $R(\mathbf{w}) = 0.5 \times \|\mathbf{w}\|_2^2$
 $\text{prox}_R(\mathbf{z}) = \mathbf{z}/2$ i.e. scaling

That is, regularize by reducing the value of each component of the vector \mathbf{z} by half

If $R(\mathbf{w})$ defines a set based constraint

$$R(\mathbf{w}) := \mathbf{w} \in \mathcal{C}$$

$$\text{prox}_R(\mathbf{z}) = \arg \min_{\mathbf{w} \in \mathcal{C}} \|\mathbf{z} - \mathbf{w}\|_2^2$$

Prox. GD becomes equivalent to projected GD



Constrained Opt. via Lagrangian

- Consider the following constrained minimization problem (using f instead of L)

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} f(\mathbf{w}), \quad \text{s.t.} \quad g(\mathbf{w}) \leq 0$$

- Note: If constraints of the form $g(\mathbf{w}) \geq 0$, use $-g(\mathbf{w}) \leq 0$
- Can handle multiple inequality and equality constraints too (will see later)
- Can transform the above into the following equivalent unconstrained problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} f(\mathbf{w}) + c(\mathbf{w})$$

$$c(\mathbf{w}) = \max_{\alpha \geq 0} \alpha g(\mathbf{w}) = \begin{cases} \infty, & \text{if } g(\mathbf{w}) > 0 \quad (\text{constraint violated}) \\ 0 & \text{if } g(\mathbf{w}) \leq 0 \quad (\text{constraint satisfied}) \end{cases}$$

- Our problem can now be written as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha g(\mathbf{w}) \right\}$$



Constrained Opt. via Lagrangian

- Therefore, we can write our original problem as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha g(\mathbf{w}) \right\} = \arg \min_{\mathbf{w}} \left\{ \arg \max_{\alpha \geq 0} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\}$$

The Lagrangian: $\mathcal{L}(\mathbf{w}, \alpha)$

- The Lagrangian is now optimized w.r.t. \mathbf{w} and α (Lagrange multiplier)
- We can defined **Primal** and **Dual** problem as

$$\begin{aligned} \hat{\mathbf{w}}_P &= \arg \min_{\mathbf{w}} \left\{ \arg \max_{\alpha \geq 0} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\} && \text{(primal problem)} \\ \hat{\mathbf{w}}_D &= \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w}} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\} && \text{(dual problem)} \end{aligned}$$

Both equal if $f(\mathbf{w})$ and the set $g(\mathbf{w}) \leq 0$ are convex

$$\alpha_D g(\hat{\mathbf{w}}_D) = 0$$

complimentary slackness/Karush-Kuhn-Tucker (KKT) condition



Constrained Opt. with Multiple Constraints

- We can also have multiple inequality and equality constraints

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} f(\mathbf{w}) \\ \text{s.t.} \quad &g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, K \\ &h_j(\mathbf{w}) = 0, \quad j = 1, \dots, L \end{aligned}$$

- Introduce Lagrange multipliers $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_K]$ and $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_L]$
- The Lagrangian based primal and dual problems will be

$$\begin{aligned} \hat{\mathbf{w}}_P &= \arg \min_{\mathbf{w}} \left\{ \arg \max_{\alpha \geq 0, \beta} \left\{ f(\mathbf{w}) + \sum_{i=1}^K \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^L \beta_j h_j(\mathbf{w}) \right\} \right\} \\ \hat{\mathbf{w}}_D &= \arg \max_{\alpha \geq 0, \beta} \left\{ \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \sum_{i=1}^K \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^L \beta_j h_j(\mathbf{w}) \right\} \right\} \end{aligned}$$



Some other useful optimization methods



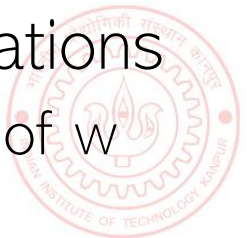
Co-ordinate Descent (CD)

- Standard gradient descent update for : $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$
- CD: In each iter, update only one entry (co-ordinate) of \mathbf{w} . Keep all others fixed

$$w_d^{(t+1)} = w_d^{(t)} - \eta_t g_d^{(t)} \quad d \in \{1, 2, \dots, D\}$$

$g_d = \nabla_{w_d} L(\mathbf{w})$ — partial derivative w.r.t. the d^{th} element of vector \mathbf{w} (or the d^{th} element of the gradient vector \mathbf{g})

- Cost of each update is now independent of D
- In each iter, can choose co-ordinate to update **unif. randomly** or in **cyclic order**
- Instead of updating a single co-ord, can also update “blocks” of co-ordinates
 - Called **Block co-ordinate descent** (BCD)
- To avoid $O(D)$ cost of gradient computation, can cache previous computations
 - Recall that grad. computations may have terms like $\mathbf{w}^T \mathbf{x}$ — if just one co-ordinate of \mathbf{w} changes, we should avoid computing the new $\mathbf{w}^T \mathbf{x}$ ($= \sum_d w_d x_d$) from scratch



Alternating Optimization (ALT-OPT)

- Consider opt. problems with several variables, say two variables \mathbf{w}_1 and \mathbf{w}_2

$$\{\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2\} = \arg \min_{\mathbf{w}_1, \mathbf{w}_2} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2)$$

- Often, this “joint” optimization is hard/impossible to solve
- We can take an alternating optimization approach to solve such problems

ALT-OPT

- 1 Initialize one of the variables, e.g., $\mathbf{w}_2 = \mathbf{w}_2^{(0)}, t = 0$
- 2 Solve $\mathbf{w}_1^{(t+1)} = \arg \min_{\mathbf{w}_1} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2^{(t)})$ // \mathbf{w}_2 “fixed” at its most recent value $\mathbf{w}_2^{(t)}$
- 3 Solve $\mathbf{w}_2^{(t+1)} = \arg \min_{\mathbf{w}_2} \mathcal{L}(\mathbf{w}_1^{(t+1)}, \mathbf{w}_2)$ // \mathbf{w}_1 “fixed” at its most recent value $\mathbf{w}_1^{(t+1)}$
- 4 $t = t + 1$. Go to step 2 if not converged yet.

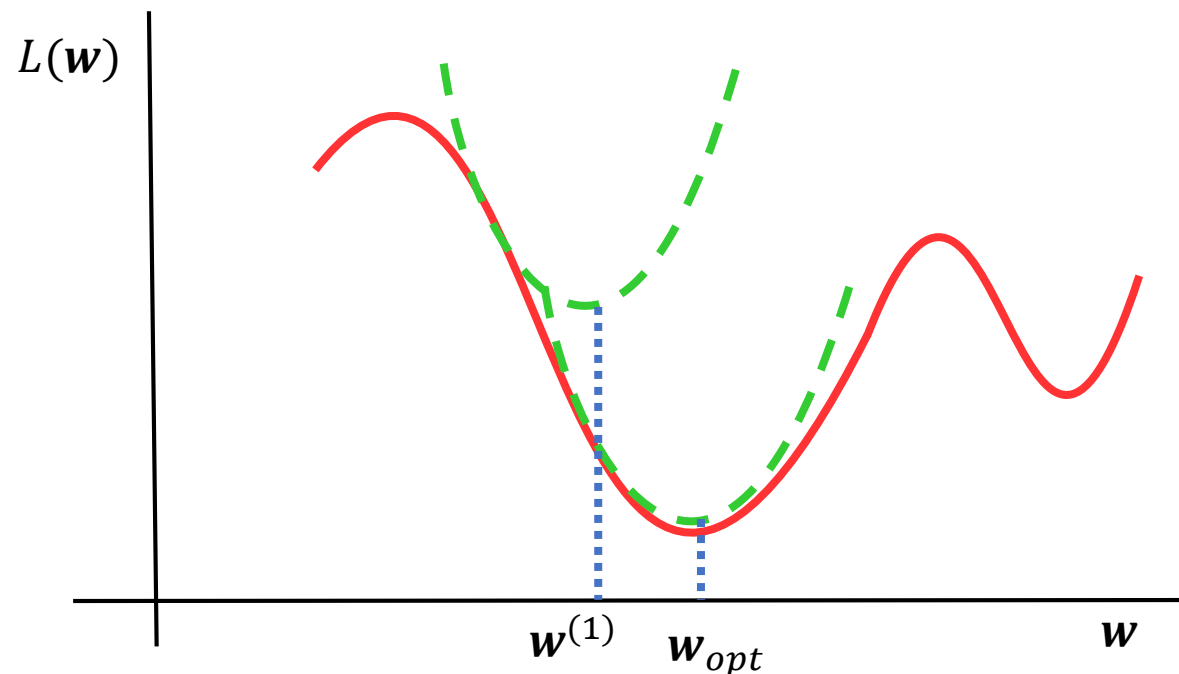
- Usually converges to a local optima. But very very useful. Will see examples later
 - Also related to the Expectation-Maximization (EM) algorithm which we will see later



Newton's Method

- Unlike GD and its variants, Newton's method uses **second-order** information (second derivative, a.k.a. the Hessian)
- At each point $\mathbf{w}^{(t)}$, minimize the quadratic (second-order) approx. of $L(\mathbf{w})$

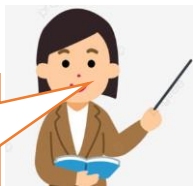
$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} [L(\mathbf{w}^{(t)}) + \nabla L(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 L(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)})]$$



Show that $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \left(\nabla^2 L(\mathbf{w}^{(t)}) \right)^{-1} \nabla L(\mathbf{w}^{(t)})$
 $= \mathbf{w}^{(t)} - (\mathbf{H}^{(t)})^{-1} \mathbf{g}^{(t)}$

Converges much faster than GD (very fast for convex functions). Also no “learning rate”. But per iteration cost is slower due to Hessian computation and inversion

Faster versions of Newton's method also exist, e.g., those based on approximating Hessian using previous gradients (see L-BFGS which is a popular method)



Coming up next

- Some practical issue in optimization for ML
- Wrapping up the discussion of optimization techniques
- Probabilistic models for ML

