

# Optimization Techniques for ML (2)

Piyush Rai

Introduction to Machine Learning (CS771A)

August 28, 2018

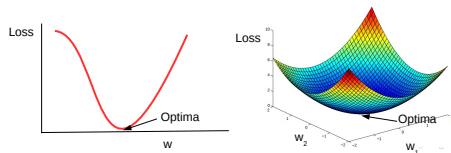


# Recap: Convex and Non-Convex Function

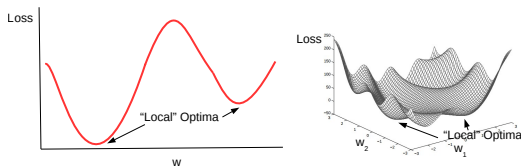
- Most ML problems boil down to minimization of convex/non-convex functions, e.g.,

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \ell_n(\mathbf{w}) + R(\mathbf{w})$$

- Convex functions have a unique minima

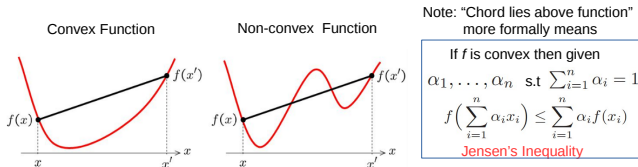


- Non-convex function have several local minima

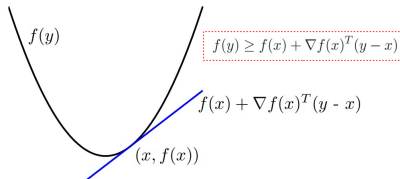


# Recap: Convex Functions

- A function is convex if all of its chords lie above the function



- A function is convex if its graph lies above all of its tangents (above its first order Taylor expansion)

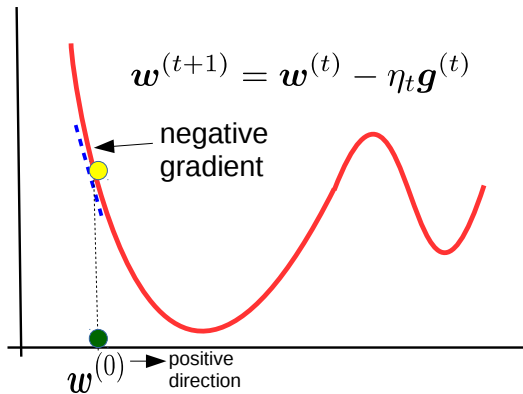


- A function is convex if its second derivative (Hessian) is positive semi-definite
- Note: If  $f$  is convex then  $-f$  is a **concave** function



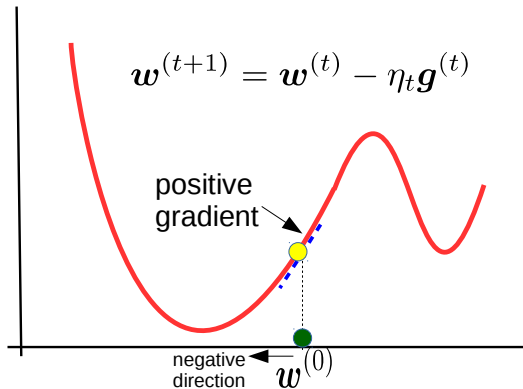
# Recap: Gradient Descent

- A very simple, **first-order method** for optimizing any differentiable function (convex/non-convex)
- Uses only the gradient  $\mathbf{g} = \nabla \mathcal{L}(\mathbf{w})$  of the function
- Basic idea: Start at some location  $\mathbf{w}^{(0)}$  and move in the **opposite direction** of the gradient



# Recap: Gradient Descent

- A very simple, **first-order method** for optimizing any differentiable function (convex/non-convex)
- Uses only the gradient  $\mathbf{g} = \nabla \mathcal{L}(\mathbf{w})$  of the function
- Basic idea: Start at some location  $\mathbf{w}^{(0)}$  and move in the **opposite direction** of the gradient



# Recap: Gradient Descent

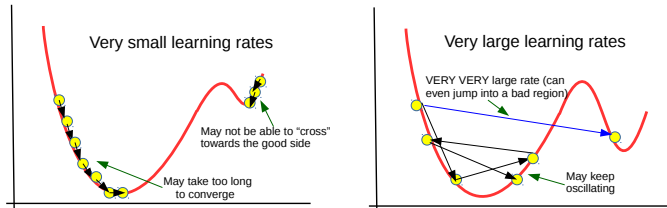
## Gradient Descent

1. Initialize  $w$  as  $w^{(0)}$
2. Update  $w$  as follows
$$w^{(t+1)} = w^{(t)} - \eta_t g^{(t)}$$
3. Repeat until convergence



# Recap: Gradient Descent

- The learning rate  $\eta_t$  is important
- Very small learning rates may result in very slow convergence
- Very large learning rates may lead to oscillatory behavior or result in a bad local optima



- Many ways to set the learning rate, e.g.,
  - Constant (if properly set, can still show good convergence behavior)
  - Decreasing with  $t$  (e.g.  $1/t$ ,  $1/\sqrt{t}$ , etc.)
  - Use **adaptive learning rates** (e.g., using methods such as **Adagrad**, **Adam**)



# Recap: Stochastic Gradient Descent

- Gradient computation in standard GD may be expensive when  $N$  is large

$$\mathbf{g} = \nabla_{\mathbf{w}} \left[ \frac{1}{N} \sum_{n=1}^N \ell_n(\mathbf{w}) \right] = \frac{1}{N} \sum_{n=1}^N \mathbf{g}_n \quad (\text{ignoring regularizer } R(\mathbf{w}))$$

- **Stochastic Gradient Descent** (SGD) approximates  $\mathbf{g}$  using a single data point
- In iteration  $t$ , SGD picks a uniformly random  $i \in \{1, \dots, N\}$  and approximate  $\mathbf{g}$  as

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

## Stochastic Gradient Descent

1. Initialize  $\mathbf{w}$  as  $\mathbf{w}^{(0)}$
2. Pick a random  $i \in \{1, \dots, N\}$ . Update  $\mathbf{w}$  as follows

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}_i^{(t)}$$

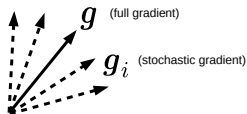
3. Repeat until convergence





# Recap: Mini-batch SGD

- In each iteration, SGD uses a single randomly chosen  $i \in \{1, \dots, N\}$  to approximate  $\mathbf{g}$
- This results in a large variance in  $\mathbf{g}_i$



- We can instead use  $B > 1$  uniformly randomly chosen points with indices  $i_1, \dots, i_B \in \{1, \dots, N\}$
- This is the idea behind mini-batch SGD. The approximated gradient in this case would be

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

- The basic intuition: Averaging helps in variance reduction!
- The algorithm is same as SGD except we will now using these mini-batch gradients at each step



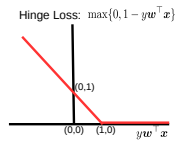
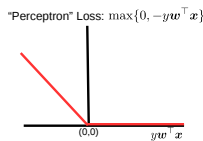
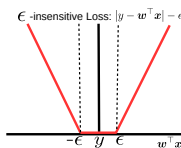
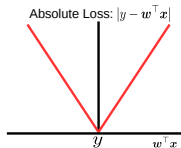
# Plan for today

- Optimization of functions that are NOT differentiable
- Optimization with constraints on the variables
- Optimizing w.r.t. several variables with one at a time
  - Co-ordinate descent
  - Alternating optimization
- Second-order methods for optimization

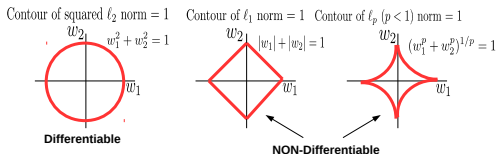


# Optimizing Non-differentiable Functions

- Many ML problems require minimizing non-differentiable functions
- Some common examples
  - Absolute,  $\epsilon$ -insensitive loss in regression, several **classification loss functions** (we will see shortly)



- Regression/classification loss functions with  $\ell_1$  or  $\ell_p$  ( $p < 1$ ) regularization



- Can't apply standard GD or SGD since gradient isn't defined at points of non-differentiability



# Interlude: Loss Functions for Classification

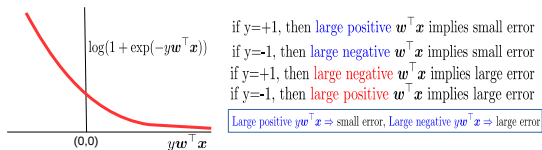
- In regression (assuming linear model  $\hat{y} = \mathbf{w}^\top \mathbf{x}$ ), some common loss functions are

$$\ell(y, \hat{y}) = (y - \mathbf{w}^\top \mathbf{x})^2 \quad \text{or} \quad \ell(y, \hat{y}) = |y - \mathbf{w}^\top \mathbf{x}|$$

- We typically look at the difference between true  $y$  and model's prediction  $\mathbf{w}^\top \mathbf{x}$
- How to formally define loss functions for **classification**?
- We have already looked at the loss function for logistic regression (assuming  $y \in \{-1, +1\}$ )

$$\ell(y, \hat{y}) = \log(1 + \exp(-y\mathbf{w}^\top \mathbf{x}))$$

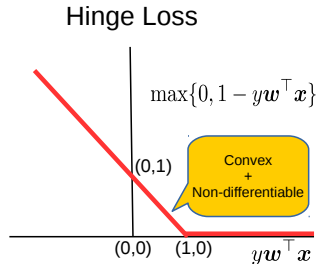
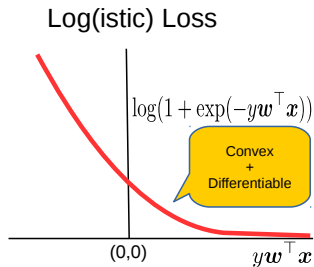
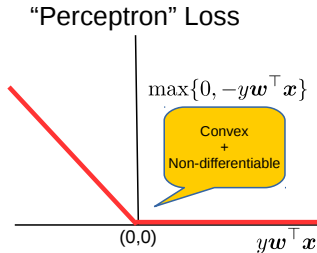
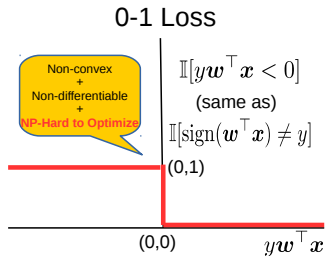
- Why does the above make sense? Well, it is large for large misclassifications, small otherwise



- Are there other loss functions for classification? Yes, several.

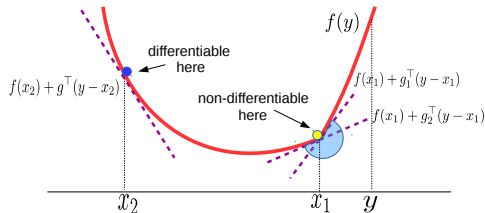


# Interlude: Some Loss Functions for (Binary) Classification



# Optimizing Non-differentiable Functions

- Even though gradients are not defined for non-diff. functions, we can work with **subgradients**



- For a function  $f(\mathbf{x})$ , its subgradient at  $\mathbf{x}$  is any vector  $\mathbf{g}$  s.t.  $\forall \mathbf{y}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x})$$

- A non-differentiable function can have **several subgradients** at the point of non-differentiability
- Set of all subgradients of a function  $f$  at point  $\mathbf{x}$  is called the **subdifferential** denoted as  $\partial f(\mathbf{x})$

$$\partial f(\mathbf{x}) = \{\mathbf{g} : f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y}\}$$



# Subgradient Descent: An Example

- Consider linear regression but with  $\ell_1$  norm on  $\mathbf{w}$  (recall:  $\ell_1$  norm promotes a sparse  $\mathbf{w}$ )

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_1$$

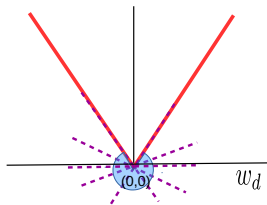
- The squared error term is differentiable but the norm  $\|\mathbf{w}\|_1$  is NOT at  $w_d = 0$
- We can use subgradients of  $\|\mathbf{w}\|_1$  in this case

$$\mathbf{g} = 2 \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n + \lambda \mathbf{t}$$

- Here  $\mathbf{t}$  is a vector s.t.

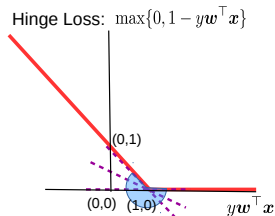
$$t_d = \begin{cases} -1, & \text{for } w_d < 0 \\ [-1, +1] & \text{for } w_d = 0 \\ +1 & \text{for } w_d > 0 \end{cases}$$

- If we take  $t_d = 0$  at  $w_d = 0$  then  $t_d = \text{sign}(w_d)$



# Subgradient Descent: Another Example

- Consider binary classification with hinge loss (used in SVM - will see later), assume  $\ell_2$  regularizer



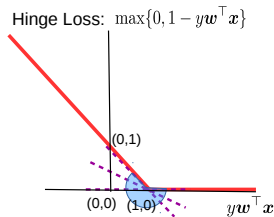
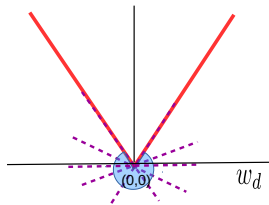
- In this case loss (hinge) non-differentiable, regularizer differentiable
- Subgradient  $\mathbf{t}$  of the hinge loss term will be

$$\mathbf{t} = \begin{cases} 0, & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n > 1 \\ -y_n \mathbf{x}_n & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n < 1 \\ ky_n \mathbf{x}_n & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n = 1 \end{cases} \quad (\text{where } k \in [-1, 0])$$





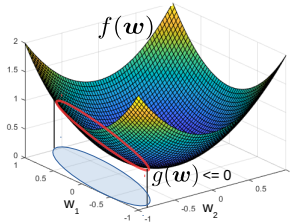
# Subgradient Descent: Summary



- Not really that different from standard GD
- Only difference is that we use subgradients where function is non-differentiable
- In practice, it is like pretending that the function is differentiable everywhere



# Constrained Optimization



- 1: Lagrangian based optimization
- 2: Projected gradient descent



# Constrained Optimization: Lagrangian Approach

- Consider optimizing some function  $f(\mathbf{w})$  subject to an **inequality constraint** on  $\mathbf{w}$

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

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

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

where we have defined  $c(\mathbf{w})$  as

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

- We can equivalently write the problem as

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



# Constrained Optimization: Lagrangian Approach

- So we could write the 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 function  $\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha g(\mathbf{w})$  called the **Lagrangian**, optimized w.r.t.  $\mathbf{w}$  and  $\alpha$
- $\alpha$  is known as the **Lagrange multiplier**
- Primal and Dual problems

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

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

- Note:  $\hat{\mathbf{w}}_P = \hat{\mathbf{w}}_D$  in some nice cases (e.g., when  $f(\mathbf{w})$  and constraint set  $g(\mathbf{w}) \leq 0$  are convex)
- For dual solution,  $\alpha_D g(\hat{\mathbf{w}}_D) = 0$  (complimentary slackness/Karush-Kuhn-Tucker (KKT) condition)

# Constrained Optimization: Lagrangian 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  $\alpha = (\alpha_1, \dots, \alpha_K) \geq 0$  and  $\beta = (\beta_1, \dots, \beta_L)$
- The Lagrangian based primal and dual problems will be

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



# Lagrangian based Optimization: An Example

- Consider the generative classification model with  $K$  classes
- Suppose we want to estimate the parameters of class-marginal  $p(y)$

$$p(y|\boldsymbol{\pi}) = \text{multinoulli}(\pi_1, \pi_2, \dots, \pi_K) = \prod_{k=1}^K \pi_k^{\mathbb{I}[y=k]}, \quad \text{s.t.} \quad \sum_{k=1}^K \pi_k = 1$$

- Given  $N$  observations  $\{\mathbf{x}_n, y_n\}_{n=1}^N$ , the negative log-likelihood for class marginal

$$f(\boldsymbol{\pi}) = - \sum_{n=1}^N \log p(y_n|\boldsymbol{\pi})$$

- We have an equality constraint  $\sum_{k=1}^K \pi_k - 1 = 0$
- The Lagrangian for this problem will be

$$\mathcal{L}(\boldsymbol{\pi}, \beta) = f(\boldsymbol{\pi}) + \beta \left( \sum_{k=1}^K \pi_k - 1 \right)$$

- **Exercise:** Solve  $\arg \max_{\beta} \arg \min_{\boldsymbol{\pi}} \mathcal{L}(\boldsymbol{\pi}, \beta)$  and show that  $\pi_k = N_k/N$

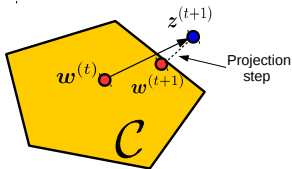


# Projected Gradient Descent

- Suppose our problem requires the parameters to lie within a set  $\mathcal{C}$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}), \quad \text{subject to} \quad \mathbf{w} \in \mathcal{C}$$

- Projected GD is very similar to GD with an extra **projection step**

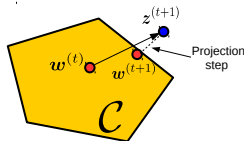


- Each step of projected GD works as follows
  - Do the usual GD update:  $\mathbf{z}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$
  - Check  $\mathbf{z}^{(t+1)}$  for the constraints
    - If  $\mathbf{z}^{(t+1)} \in \mathcal{C}$ ,  $\mathbf{w}^{(t+1)} = \mathbf{z}^{(t+1)}$
    - If  $\mathbf{z}^{(t+1)} \notin \mathcal{C}$ , project on the constraint set:  $\mathbf{w}^{(t+1)} = \underbrace{\Pi_{\mathcal{C}}[\mathbf{z}^{(t+1)}]}_{\text{projection}}$



# Projected GD: How to Project?

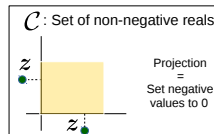
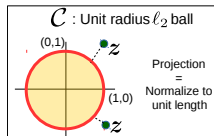
- The projection itself is an optimization problem



- Given  $\mathbf{z}$ , we find the “closest” point (e.g., in Euclidean sense)  $\mathbf{w}$  in the set as follows

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

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



- For some other sets  $\mathcal{C}$ , the projection step may be a bit more involved



# Co-ordinate Descent (CD)

- Standard GD update for  $\mathbf{w} \in \mathbb{R}^D$  at each step

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

- CD: Each step update one component (co-ordinate) at a time, keeping all others fixed

$$w_d^{(t+1)} = w_d^{(t)} - \eta_t g_d^{(t)}$$

- Cost of each update is now independent of  $D$
- How to pick which co-ordinate to update?
  - Can be chosen in random order (stochastic CD)
  - Can be chosen in cyclic order
- Note: Can also update “blocks” of co-ordinates (called Block co-ordinate descent)
- Should cache previous computations (e.g.,  $\mathbf{w}^\top \mathbf{x}$ ) to avoid  $\mathcal{O}(D)$  cost in gradient computation



# Alternating Optimization

- Many optimization problems consist of several variables. Very common in ML.
- For simplicity, suppose we want to optimize a function of 2 variables  $\mathbf{w}_1 \in \mathbb{R}^D$  and  $\mathbf{w}_2 \in \mathbb{R}^D$

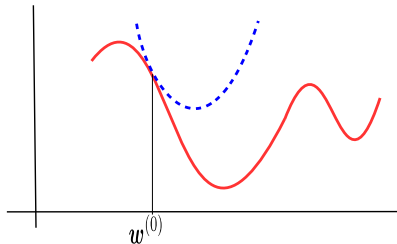
$$\{\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2\} = \arg \min_{\mathbf{w}_1, \mathbf{w}_2} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2)$$

- Jointly optimizing w.r.t.  $\mathbf{w}_1$  and  $\mathbf{w}_2$  may be hard (e.g., if their values depend on each other)
- Often, knowing value of one may make optimization w.r.t. other easy (sometimes even **closed form**)
- We can therefore follow an **alternating scheme** to optimize w.r.t.  $\mathbf{w}_1$  and  $\mathbf{w}_2$ 
  - Initialize one of the variables, e.g.,  $\mathbf{w}_2 = \mathbf{w}_2^{(0)}, t = 0$
  - Solve  $\mathbf{w}_1^{(t+1)} = \arg \max_{\mathbf{w}_1} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2^{(t)})$
  - Solve  $\mathbf{w}_2^{(t+1)} = \arg \max_{\mathbf{w}_2} \mathcal{L}(\mathbf{w}_1^{(t+1)}, \mathbf{w}_2)$
  - $t = t + 1$ . Repeat until convergence
- Usually converges to a **local optima** of  $\mathcal{L}(\mathbf{w}_1, \mathbf{w}_2)$ . Also connections to **EM** (will see later)
- Extends to more than 2 variables as well (and not just to vectors). CD is a special case.



# Second-Order Methods: Newton's Method

- GD and variants only use first-order information (the gradient)
- Second-order information often tells us a lot more about the function's shape, curvature, etc.
- Newton's method is one such method that uses second-order information
- At each point, approximate the function by its quadratic approx. and minimize it

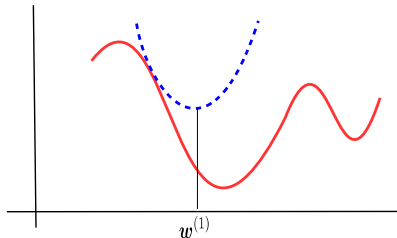


- Doesn't rely on gradient to choose  $\mathbf{w}^{(t+1)}$
- Instead, each step directly jumps to the minima of quadratic approximation



# Second-Order Methods: Newton's Method

- GD and variants only use first-order information (the gradient)
- Second-order information often tells us a lot more about the function's shape, curvature, etc.
- Newton's method is one such method that uses second-order information
- At each point, approximate the function by its quadratic approx. and minimize it

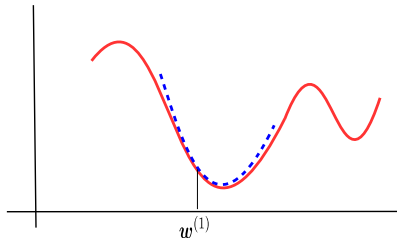


- Doesn't rely on gradient to choose  $\mathbf{w}^{(t+1)}$
- Instead, each step directly jumps to the minima of quadratic approximation



# Second-Order Methods: Newton's Method

- GD and variants only use first-order information (the gradient)
- Second-order information often tells us a lot more about the function's shape, curvature, etc.
- Newton's method is one such method that uses second-order information
- At each point, approximate the function by its quadratic approx. and minimize it

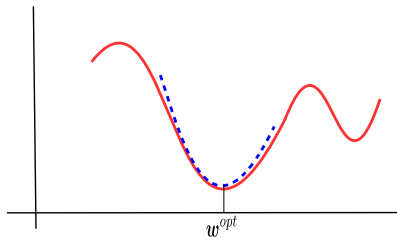


- Doesn't rely on gradient to choose  $\mathbf{w}^{(t+1)}$
- Instead, each step directly jumps to the minima of quadratic approximation



# Second-Order Methods: Newton's Method

- GD and variants only use first-order information (the gradient)
- Second-order information often tells us a lot more about the function's shape, curvature, etc.
- Newton's method is one such method that uses second-order information
- At each point, approximate the function by its quadratic approx. and minimize it



- Doesn't rely on gradient to choose  $\mathbf{w}^{(t+1)}$
- Instead, each step directly jumps to the minima of quadratic approximation



# Second-Order Methods: Newton's Method

- The quadratic (Taylor) approximation of  $f(\mathbf{w})$  at  $\mathbf{w}^{(t)}$  is given by

$$\tilde{f}(\mathbf{w}) = f(\mathbf{w}^{(t)}) + \nabla f(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 f(\mathbf{w}^{(t)})(\mathbf{w} - \mathbf{w}^{(t)})$$

- The minimizer of this quadratic approximation is (exercise: verify)

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \tilde{f}(\mathbf{w}) = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$$

- This is the update used in Newton's method (a second order method since it uses the Hessian)

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$$

- Look, Ma! No learning rate! :-)
- Very fast if  $f(\mathbf{w})$  is convex. But expensive due to Hessian computation/inversion.
- Many ways to approximate the Hessian (e.g., using previous gradients); also look at L-BFGS etc.



# Summary

- Gradient methods are simple to understand and implement
- More sophisticated optimization methods often use gradient methods
  - **Backpropagation algorithm** used in deep neural nets is **GD + chain rule** of differentiation
- Use **subgradient** methods if function not differentiable
- Constrained optimization require methods such as **Lagrangian** or **projected gradient**
- **Second order methods** such as Newton's method are much faster but computationally expensive
- But computing all this gradient related stuff looks scary to me. Any help?
  - Don't worry. **Automatic Differentiation** (AD) methods available now
  - AD only requires specifying the loss function (especially useful for deep neural nets)
  - Many packages such as Tensorflow, PyTorch, etc. provide AD support
  - But having a good understanding of optimization is still helpful

