

Optimization for ML

CS771: Introduction to Machine Learning

Plan today

- Basics of optimization of functions
 - Global and local optima
 - First derivative (gradient) and second derivative (Hessian)
- Convex and non-convex functions
- Popular optimization methods
 - Gradient descent
 - Stochastic gradient descent



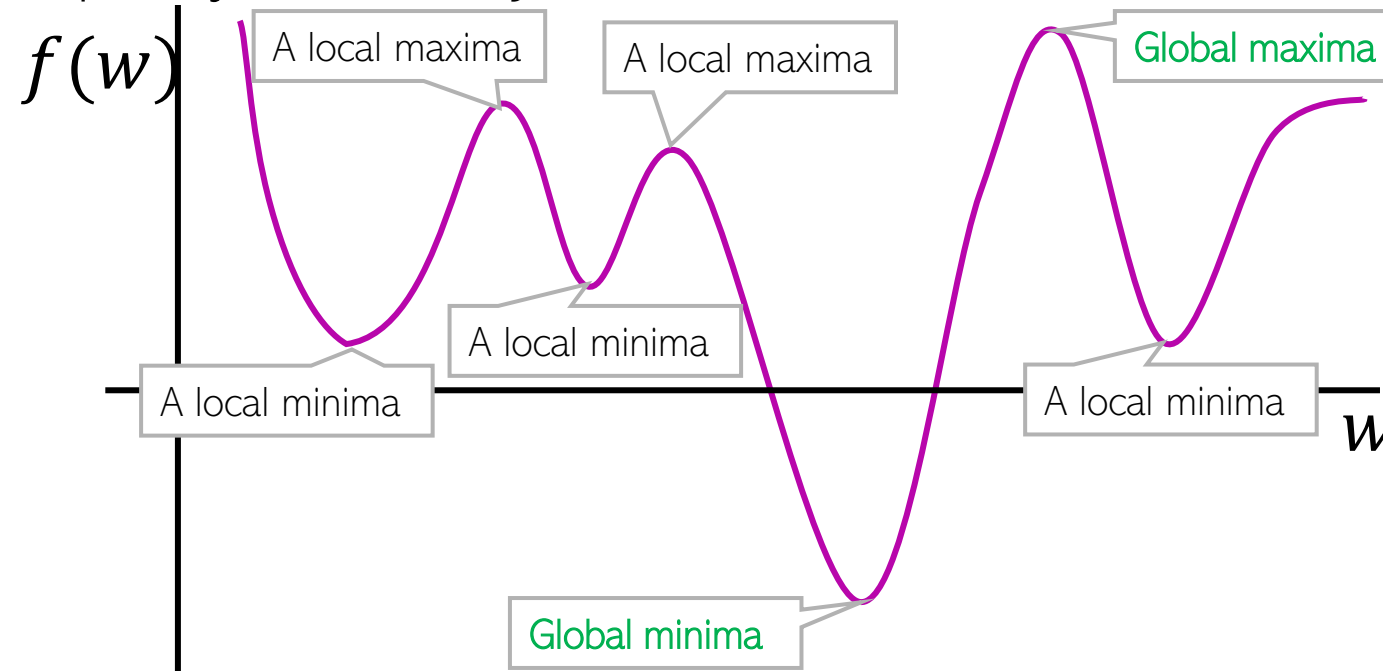
Functions and their optima

The objective function of the ML problem we are solving (e.g., squared loss for regression)

Assume **unconstrained** for now, i.e., just a real-valued number/vector

3

- Many ML problems require us to optimize a function f of some variable(s) w
- For simplicity, assume f is a scalar-valued function of a scalar w ($f: \mathbb{R} \rightarrow \mathbb{R}$)



Usually interested in global optima but often want to find local optima, too

For deep learning models, often the local optima are what we can find (and they usually suffice) – more later

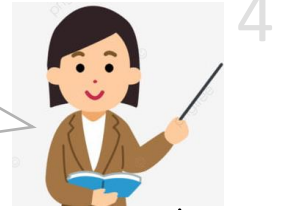
Will see what these are later

- Any function has one/more **optima** (maxima, minima), and maybe **saddle points**
- Finding the optima or saddles requires derivatives/gradients of the function



Derivatives

Will sometimes use $f'(w)$ to denote the derivative

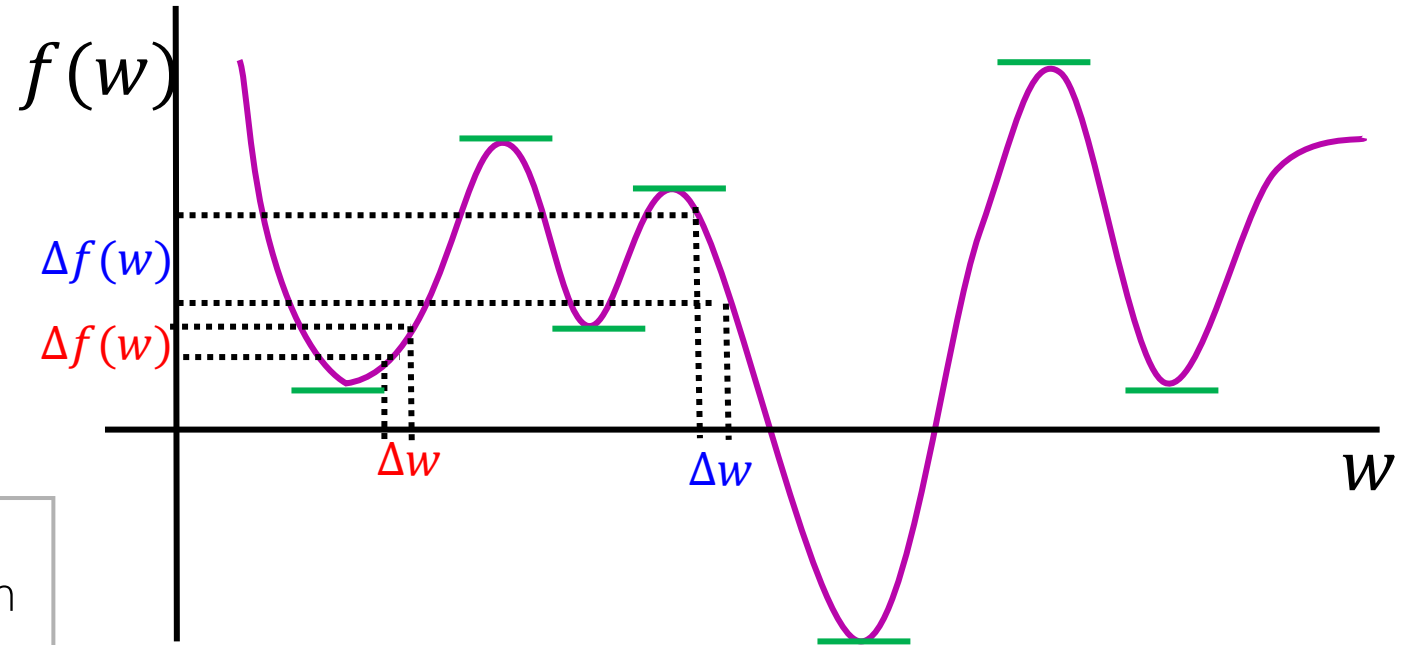


- Magnitude of derivative at a point is the rate of change of the func at that point

$$\frac{df(w)}{dw} = \lim_{\Delta w \rightarrow 0} \frac{\Delta f(w)}{\Delta w}$$

Sign is also important: Positive derivative means f is **increasing** at w if we increase the value of w by a very small amount; negative derivative means it is **decreasing**

Understanding how f changes its value as we change w is helpful to understand optimization (minimization/maximization) algorithms

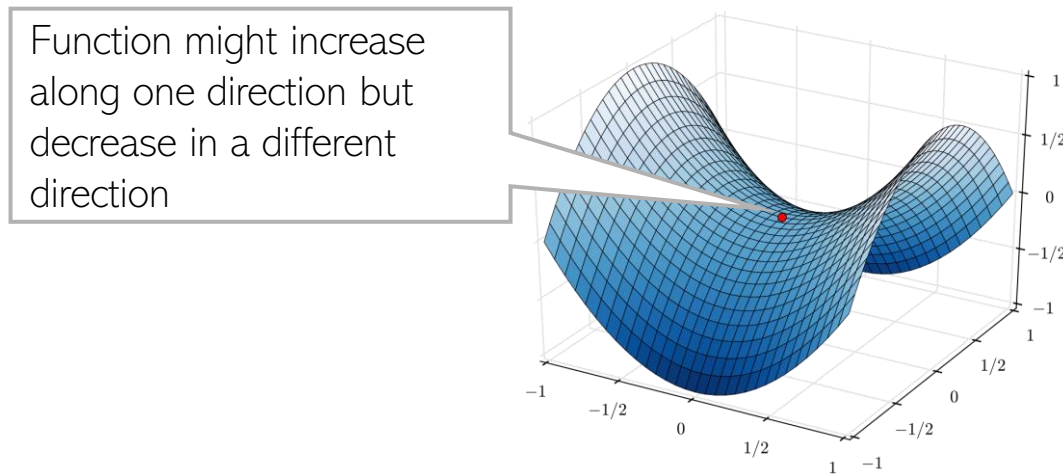


- Derivative becomes zero at **stationary points** (**optima** or **saddle points**)
 - The function becomes **“flat”** ($\Delta f(w) = 0$ if we change w by a very little at such points)
 - These are the points where the function has its maxima/minima (unless they are saddles)

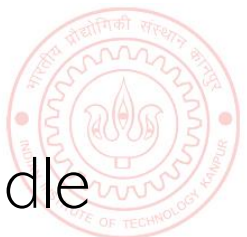


Saddle Points

- Points where derivative is zero but are neither minima nor maxima



- Saddle points are very common for loss functions of deep learning models
 - Need to be handled carefully during optimization
- Second or higher derivative may help identify if a stationary point is a saddle



Rules of Derivatives

Some basic rules of taking derivatives

- Sum Rule: $(f(w) + g(w))' = f'(w) + g'(w)$
- Scaling Rule: $(a \cdot f(w))' = a \cdot f'(w)$ if a is not a function of w
- Product Rule: $(f(w) \cdot g(w))' = f'(w) \cdot g(w) + g'(w) \cdot f(w)$
- Quotient Rule: $(f(w)/g(w))' = (f'(w) \cdot g(w) - g'(w)f(w))/(g(w))^2$
- Chain Rule: $(f(g(w)))' \stackrel{\text{def}}{=} (f \circ g)'(w) = f'(g(w)) \cdot g'(w)$

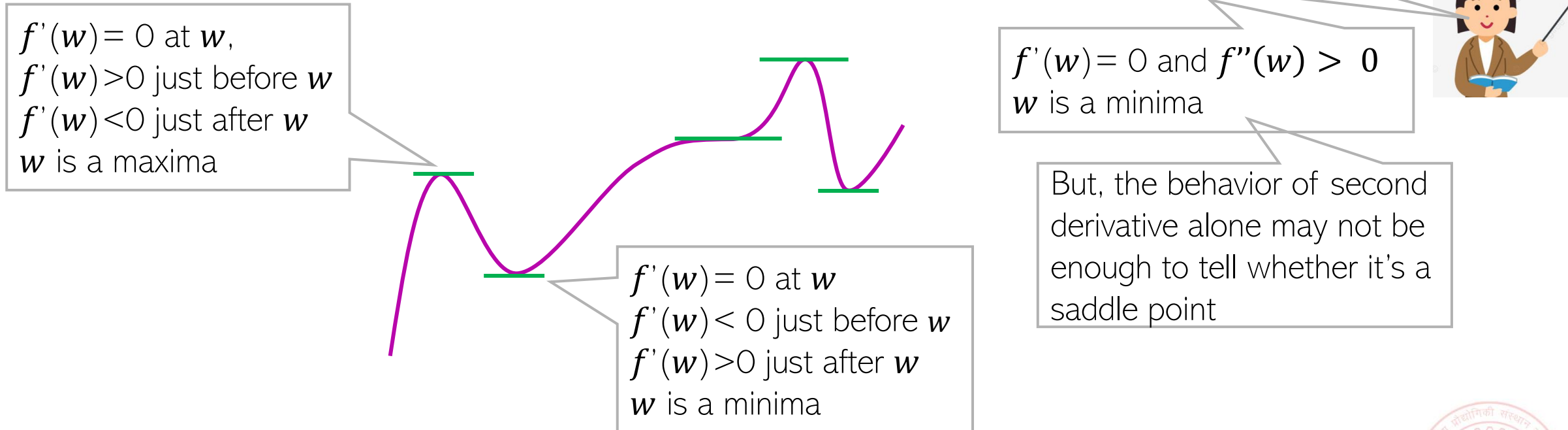


We already used some of these (sum, scaling and chain) when calculating the derivative for the linear regression model



Derivatives

- How the derivative itself changes tells us about the function's optima



- The second derivative $f''(w)$ can provide this information
 - It is the rate of change of the first derivative

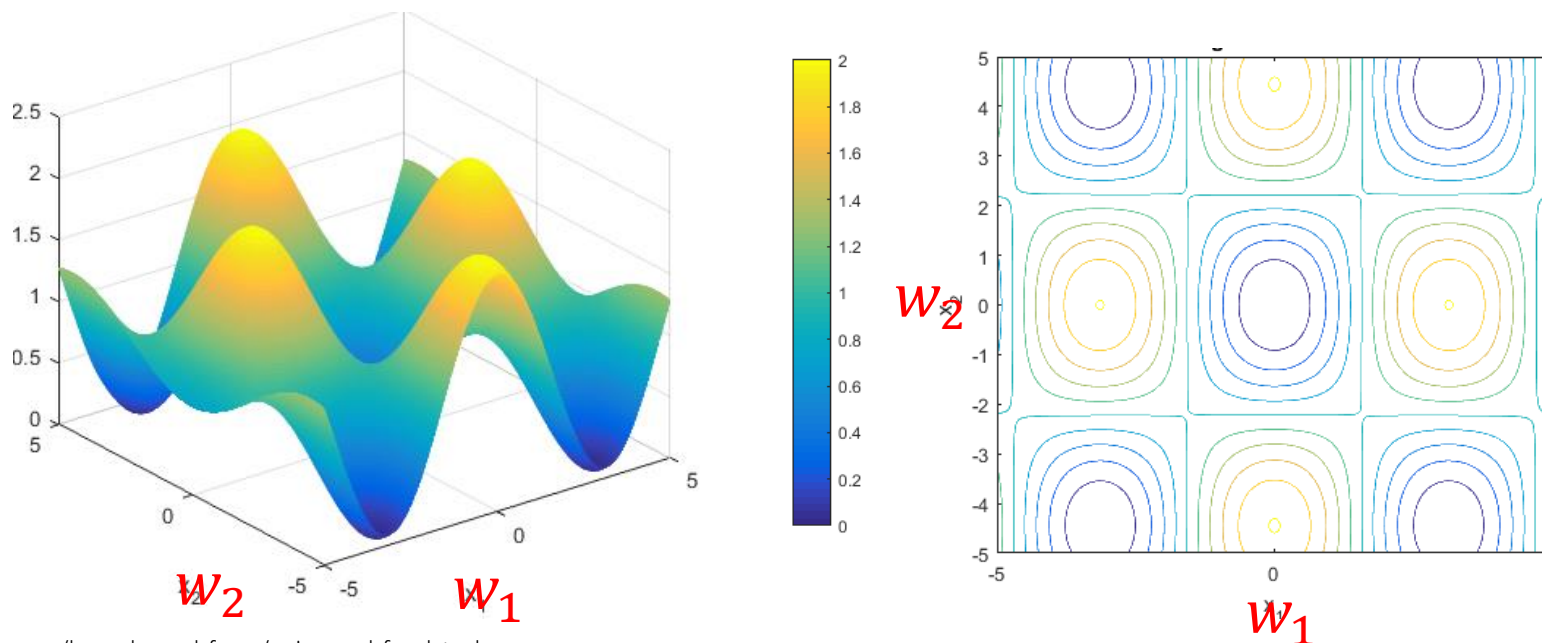


Multivariate Functions

- Most functions that we see in ML are multivariate function
- Example: Loss fn $L(\mathbf{w})$ in lin-reg was a multivar function of D -dim vector \mathbf{w}

$$L(\mathbf{w}): \mathbb{R}^D \rightarrow \mathbb{R}$$

- Here is an illustration of a function of 2 variables (4 maxima and 5 minima)



Two-dim contour plot of the function (i.e., what it looks like from the above)

Each curve here is the set of points at which the function has the same value

Derivatives of Multivariate Functions

- Can define derivative for a multivariate functions as well via the gradient
- Gradient of a function $f(\mathbf{w}): \mathbb{R}^D \rightarrow \mathbb{R}$ is a $D \times 1$ vector of partial derivatives

$$\nabla f(\mathbf{w}) = \left(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_D} \right)^T$$

Each element in this gradient vector tells us how much f will change if we move a little along the corresponding (akin to one-dim case)

- Conditions for optima defined similar to one-dim case
 - Required properties for one-dim case must be satisfied for all components of \mathbf{w}
- If the function itself is vector valued, e.g., $f(\mathbf{w}): \mathbb{R}^D \rightarrow \mathbb{R}^K$ then we will have K such gradient vectors, one for each output dimension of f
- The second derivative in the multivariate case is known as the **Hessian matrix**



The Hessian

- For a multivar scalar valued function $f(\mathbf{w}): \mathbb{R}^D \rightarrow \mathbb{R}$, Hessian is a $D \times D$ matrix

$$\nabla^2 f(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 w_D} \\ \frac{\partial^2 f}{\partial w_2 w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 w_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_D w_1} & \frac{\partial^2 f}{\partial w_D w_2} & \cdots & \frac{\partial^2 f}{\partial w_D^2} \end{bmatrix}$$

Gives information about the **curvature** of the function at point \mathbf{w}

Note: If the function itself is vector valued, e.g., $f(\mathbf{w}): \mathbb{R}^D \rightarrow \mathbb{R}^K$ then we will have K such $D \times D$ Hessian matrices, one for each output dimension of f

A square, symmetric $D \times D$ matrix \mathbf{M} is PSD if $\mathbf{w}^T \mathbf{M} \mathbf{w} \geq 0 \forall \mathbf{w} \in \mathbb{R}^D$
Will be NSD if $\mathbf{w}^T \mathbf{M} \mathbf{w} \leq 0 \forall \mathbf{w} \in \mathbb{R}^D$

PSD if all eigenvalues are non-negative

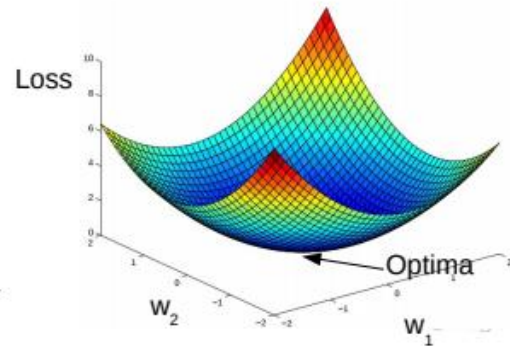
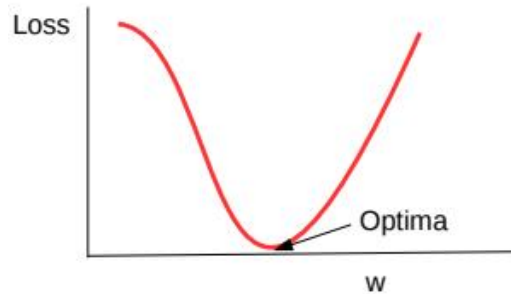


- The Hessian matrix can be used to assess the optima/saddle points
 - $\nabla f(\mathbf{w}) = 0$ and $\nabla^2 f(\mathbf{w})$ is a positive semi-definite (PSD) matrix then \mathbf{w} is a minima
 - $\nabla f(\mathbf{w}) = 0$, and $\nabla^2 f(\mathbf{w})$ is a negative semi-definite (NSD) matrix then \mathbf{w} is a maxima
 - The determinant of the Hessian is zero at saddle point (has at least one zero eig-val)



Convex and Non-Convex Functions

- A function being optimized can be either **convex** or **non-convex**
- Here are a couple of examples of convex functions

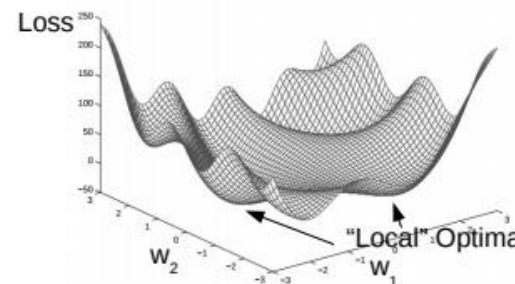
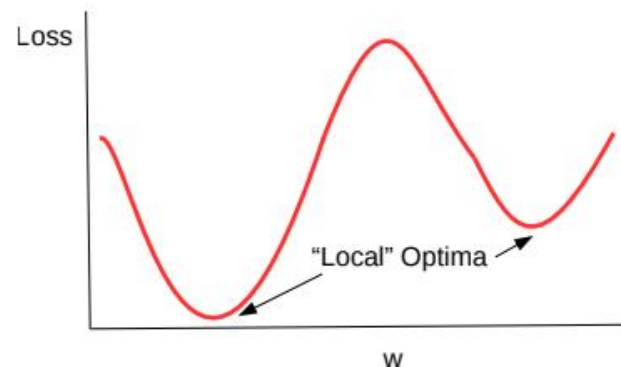


Convex functions are bowl-shaped.
They have a unique optima (minima)

Negative of a convex function is called
a **concave** function, which also has a
unique optima (maxima)



- Here are a couple of examples of non-convex functions



Non-convex functions have
multiple minima. Usually harder
to optimize as compared to
convex functions

Loss functions of most
deep learning models are
non-convex



Convex Sets

- A set S of points is a convex set, if for any two $x, y \in S$, and $0 \leq \lambda \leq 1$

z is also called a “convex combination” of two points

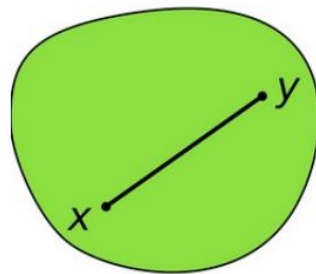
$$z = \lambda x + (1 - \lambda)y \in S$$

Can also define **convex combination** of N points x_1, x_2, \dots, x_N using non-neg $\lambda_1, \lambda_2, \dots, \lambda_N$ s.t., $\sum_{i=1}^N \lambda_i = 1$ as
 $z = \sum_{i=1}^N \lambda_i x_i$

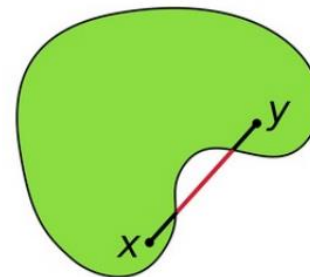


- Above means that all points on the line-segment between x and y lie within S

A Convex Set



A Non-convex Set

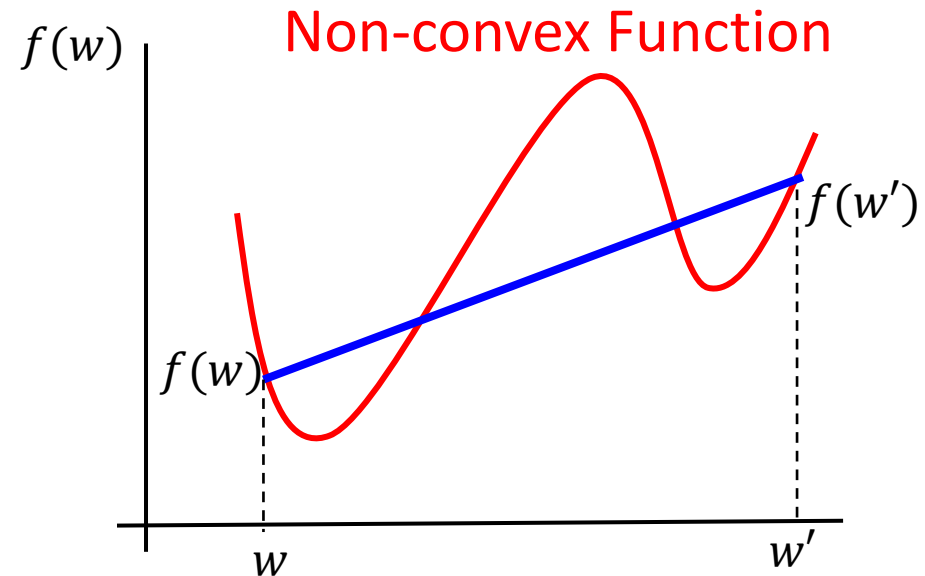
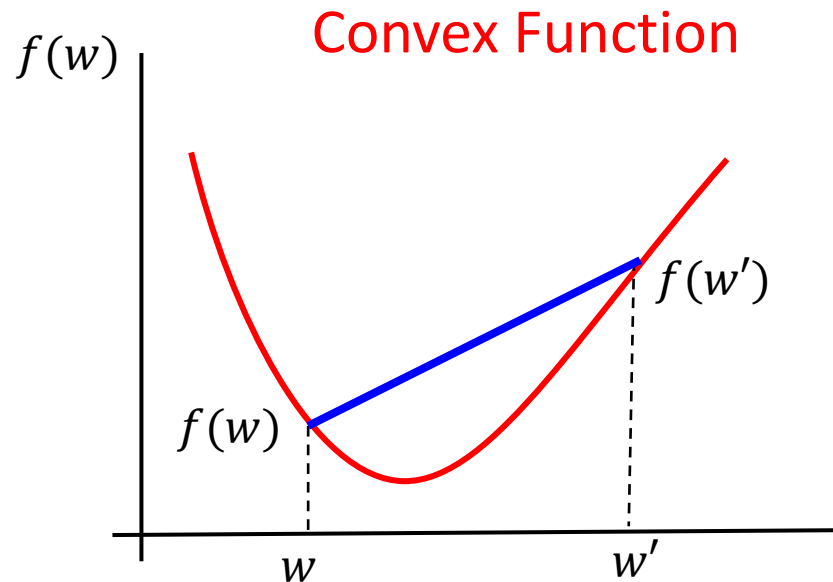


- The domain of a convex function needs to be a convex set



Checking if a function is convex

- Informally, $f(w)$ is convex if ALL of its chords lie above the function everywhere



- If $f(w)$ is convex then for any $\lambda \in (0,1)$ and all pairs of points w and w'

"Chord lies above"
condition, mathematically

$$f(\lambda w + (1 - \lambda)w') \leq \lambda f(w) + (1 - \lambda)f(w')$$

Even more general, holds
for **convex combination** of
any n points

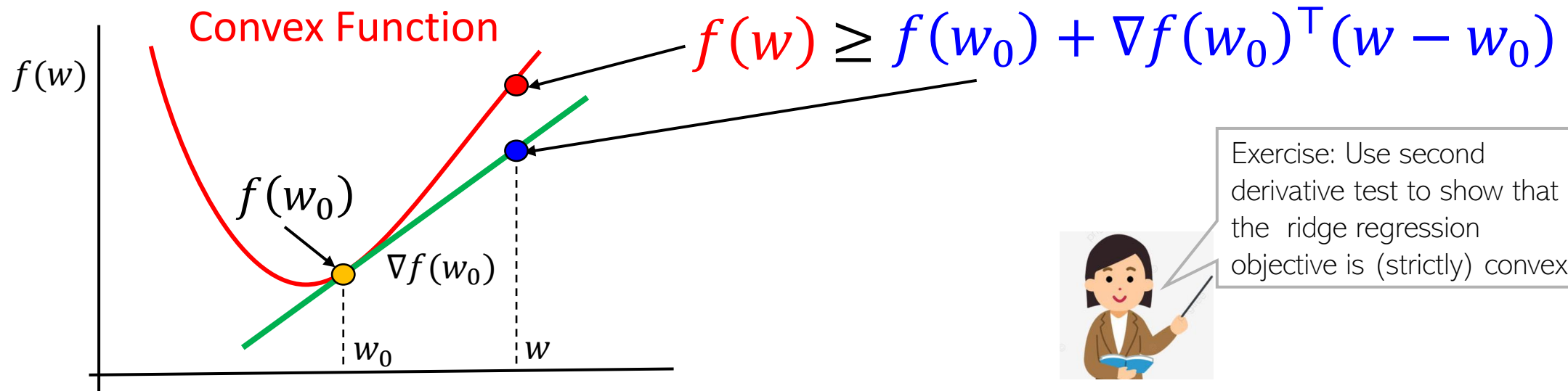
$$f\left(\sum_{i=1}^n \lambda_i w_i\right) \leq \sum_{i=1}^n \lambda_i f(w_i)$$

Jensen's Inequality



Checking for Convexity using Derivatives

- First-order convexity (graph of f must be above all the tangents)



- Second derivative a.k.a. Hessian matrix(if exists) must be positive semi-definite

$$\nabla^2 f(\mathbf{w}) \succcurlyeq \mathbf{0}$$

- Strictly convex: Above inequalities are strict (i.e., not \geq but $>$)



Some Basic Rules for Convex Functions

- All linear or affine functions ($aw + b$) are convex in w for any $a, b \in \mathbb{R}$
- $\exp(aw)$ is convex in $w \in \mathbb{R}$ for any $a \in \mathbb{R}$
- $\log(w)$ is **concave** in $w > 0$ ($-\log(w)$ is convex)
- w^a is convex in $w > 0$ for any $a \geq 1$, and concave for $0 \leq a \leq 1$
- $|w|^a$ is convex in $w \in \mathbb{R}$ for any $a \geq 1$
- All norm functions (e.g., ℓ_2 -squared norm, ℓ_1 norm) are convex
- Nonnegative weighted sum of convex functions is also convex function
- If $f(w)$ is convex then $f(aw + b)$ is also convex
- Rules regarding compositions of two functions h and g : $f(w) = h(g(w))$

h convex and nondecreasing, g convex: f convex

h convex and nondecreasing, g concave: f concave

h convex and nonincreasing, g convex: f concave

h convex and nonincreasing, g concave: f convex

Optimization Problems in ML

- The general form of an optimization problem in ML will usually be

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

$L(\mathbf{w})$ may denote the training loss,
or training loss + regularizer term

or

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

Training loss with a
constraint on \mathbf{w}

- \mathcal{C} is the constraint set that the solution must belong to, e.g.,
 - Non-negativity constraint: All entries in \mathbf{w}_{opt} must be non-negative
 - Sparsity constraint: \mathbf{w}_{opt} is a sparse vector with at most K non-zeros
- Constrained opt. probs can be converted into unconstrained opt. (will see later)
- For now, assume we have an unconstrained optimization problem

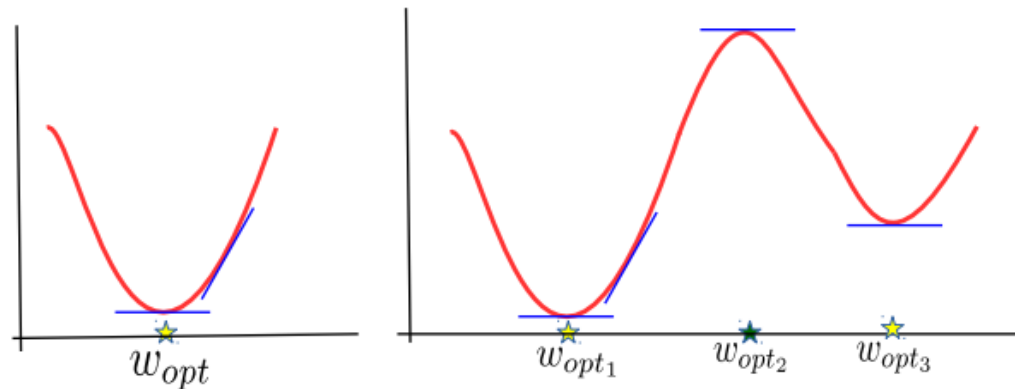


Methods for Solving Optimization Problems



Method 1: Using First-Order Optimality

- Very simple. Already used this approach for linear and ridge regression



Called “first order” since only gradient is used and gradient provides the first order info about the function being optimized



The approach works only for very simple problems where the objective is convex and there are no constraints on the values \mathbf{w} can take

- First order optimality: The gradient \mathbf{g} must be equal to zero at the optima

$$\mathbf{g} = \nabla_{\mathbf{w}}[L(\mathbf{w})] = \mathbf{0}$$

E.g., linear/ridge regression, but not for logistic/softmax regression

- Sometimes, setting $\mathbf{g} = \mathbf{0}$ and solving for \mathbf{w} gives a closed form solution
- If closed form solution is not available, the gradient vector \mathbf{g} can still be used in iterative optimization algos, like [gradient descent](#)



Method 2: Iterative Optimiz. via Gradient Descent¹⁹



Can I used this approach to solve **maximization** problems?

For max. problems we can use gradient **ascent**
 $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta_t \mathbf{g}^{(t)}$

Iterative since it requires several steps/iterations to find the optimal solution



For convex functions, GD will converge to the global minima

Good initialization needed for non-convex functions

Fact: Gradient gives the direction of **steepest change** in function's value

Will move in the direction of the gradient

Gradient Descent

- Initialize \mathbf{w} as $\mathbf{w}^{(0)}$
- For iteration $t = 0, 1, 2, \dots$ (or until convergence)
 - Calculate the gradient $\mathbf{g}^{(t)}$ using the current iterates $\mathbf{w}^{(t)}$
 - Set the learning rate η_t
 - Move in the **opposite** direction of gradient

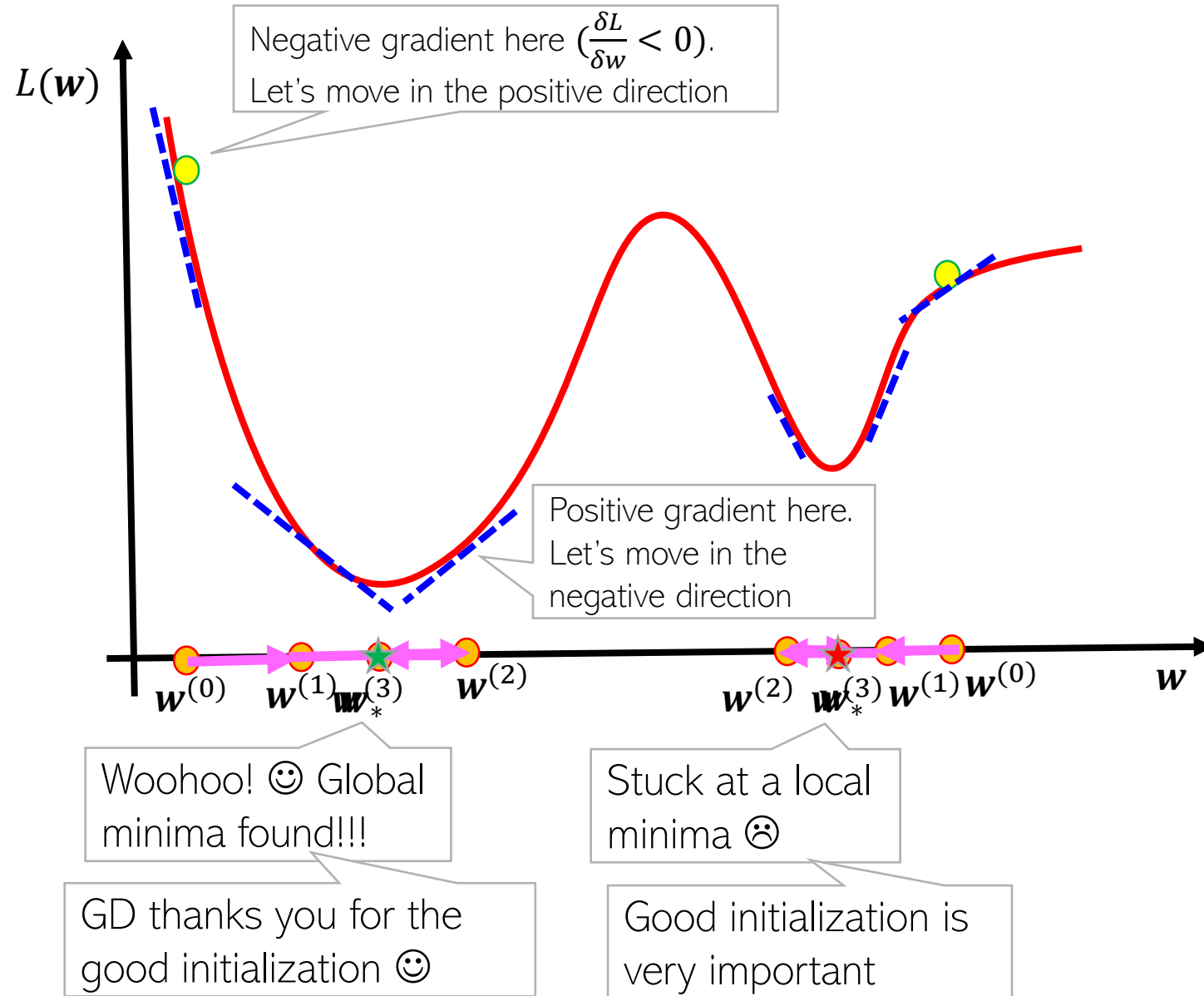
Will see the justification shortly

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

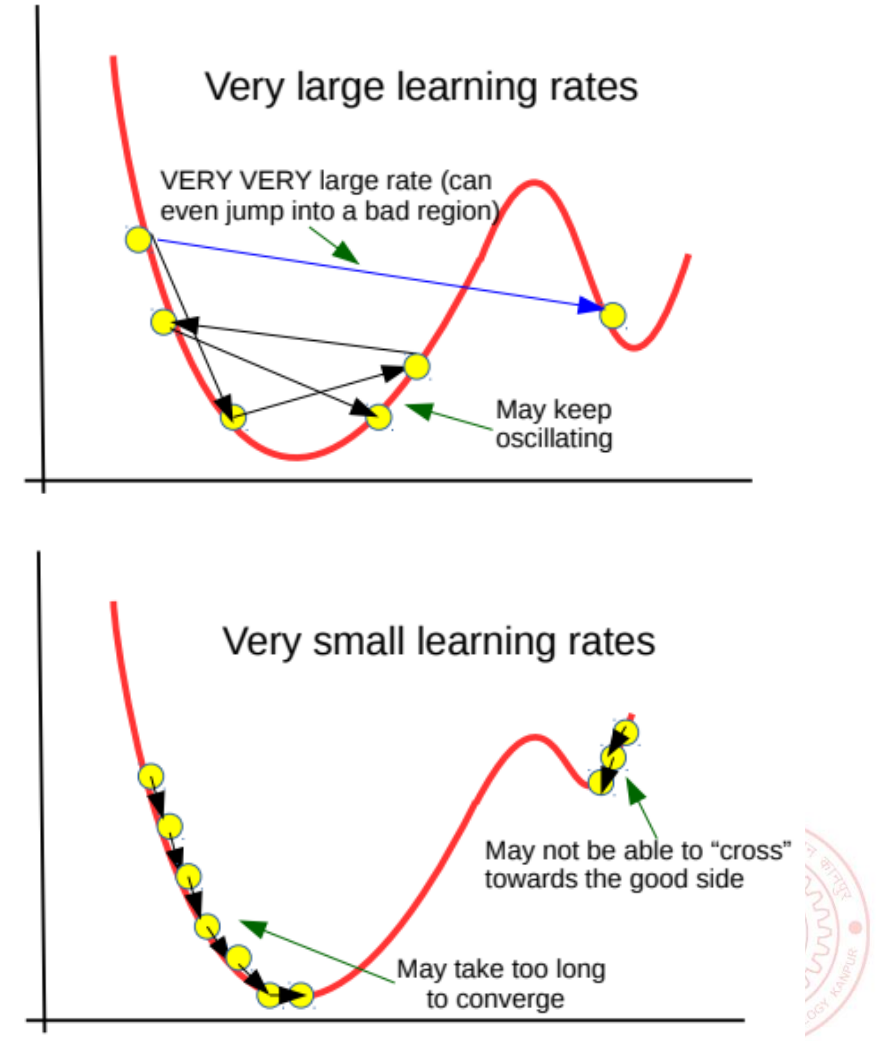
The learning rate very imp. Should be set carefully (fixed or chosen adaptively). Will discuss some strategies later



Gradient Descent: An Illustration



Learning rate is very important



Faster GD: Stochastic Gradient Descent (SGD)

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

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

- The 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 ☹

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

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

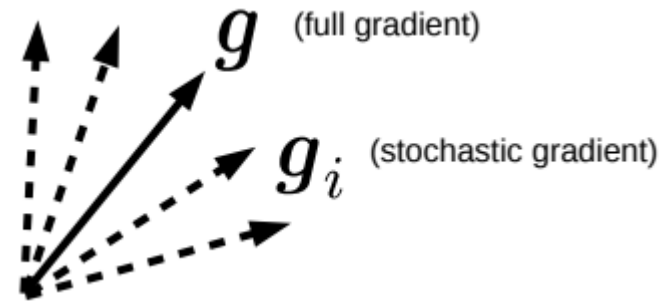
Can show that \mathbf{g}_i is an **unbiased estimate** of \mathbf{g} ,
i.e., $\mathbb{E}[\mathbf{g}_i] = \mathbf{g}$

- 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

