

03-Linear Regression (Convexity)

Convexity

Definition of convexity

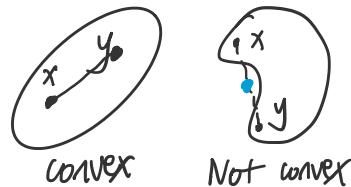
Convex set

Let X be a vector space, and $S \subseteq X$

S is a convex set \Leftrightarrow

$$\forall \mathbf{x}, \mathbf{y} \in S \quad \forall \lambda \in (0, 1)$$

$$\lambda \mathbf{x} + (1-\lambda) \mathbf{y} \in S$$



Convex function

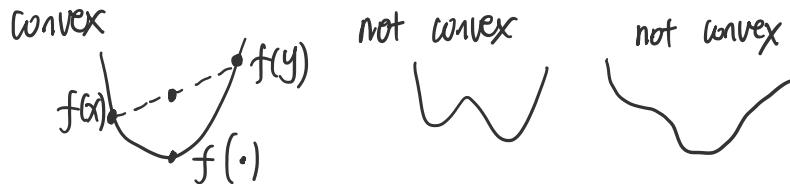
f is a convex function \Leftrightarrow

$\text{dom } f$ is a convex set [domain of f]

AND

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom } f, \quad \forall \lambda \in (0, 1)$$

$$\lambda f(\mathbf{x}) + (1-\lambda) f(\mathbf{y}) \geq f(\lambda \mathbf{x} + (1-\lambda) \mathbf{y})$$

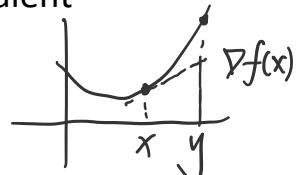


Theorem:

Let f be a twice-differentiable function on a convex open domain $\text{dom } f$.

Then, the following statements are equivalent

a. f is convex



b. (First-order condition)

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom } f, \quad f(\mathbf{y}) \geq f(\mathbf{x}) + [\nabla f(\mathbf{x})]^T \cdot (\mathbf{y} - \mathbf{x})$$

c. (Second-order condition)

$$\forall \mathbf{x} \in \text{dom } f \quad \nabla^2 f(\mathbf{x}) \succeq 0 \quad \text{positive semi-definite}$$

Calculus review

The **gradient** is the first-order partial derivatives arranged in the original shape, denoted by $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ or $\nabla_{\mathbf{x}} f(\mathbf{x})$:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}) \\ \frac{\partial}{\partial x_2} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_n} f(\mathbf{x}) \end{pmatrix} \in \mathbb{R}^n, \text{ where } \mathbf{x} \in \mathbb{R}^n$$

The **Hessian** is the second-order partial derivatives, arranged as a matrix

denoted by $\mathbf{H}(\mathbf{x})$ or $\nabla^2 f(\mathbf{x})$

$$\mathbf{H}(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \ddots & & \\ \vdots & \ddots & \ddots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}_{n \times n}$$

Proof: b⇒a; c⇒b in one-dimensional case. Others=homework

(Direction: b⇒a)

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom } f \quad \text{and} \quad \forall \lambda \in (0, 1)$$

$$\text{Def } \mathbf{z} := \lambda \mathbf{x} + (1-\lambda) \mathbf{y}$$

Applying condition b:

$$f(\mathbf{x}) \geq f(\mathbf{z}) + [\nabla f(\mathbf{z})]^T \cdot (\mathbf{x} - \mathbf{z}) \quad (1)$$

$$f(\mathbf{y}) \geq f(\mathbf{z}) + [\nabla f(\mathbf{z})]^T \cdot (\mathbf{y} - \mathbf{z}) \quad (2)$$

$$(1) \cdot \lambda: \quad \lambda f(\mathbf{x}) \geq \lambda f(\mathbf{z}) + \lambda \cdot [\nabla f(\mathbf{z})]^T \cdot (\mathbf{x} - \mathbf{z}) \quad (3)$$

$$(2) \cdot (1-\lambda): \quad (1-\lambda) f(\mathbf{y}) \geq (1-\lambda) f(\mathbf{z}) + (1-\lambda) [\nabla f(\mathbf{z})]^T \cdot (\mathbf{y} - \mathbf{z}) \quad (4)$$

(3) + (4):

$$\begin{aligned} \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y}) &\geq f(\mathbf{z}) + [\nabla f(\mathbf{z})]^T \cdot (\underbrace{\lambda \cdot (\mathbf{x}-\mathbf{z}) + (1-\lambda)(\mathbf{y}-\mathbf{z})}_{\mathbf{z}}) \\ &= f(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \quad \underbrace{\lambda \mathbf{x} + (1-\lambda)\mathbf{y} - \mathbf{z}}_{\mathbf{z}} = \lambda \mathbf{x} + (1-\lambda)\mathbf{y} - \mathbf{z} \end{aligned}$$

(Direction: $\Rightarrow b$ in one-dimensional case)

By Taylor's Theorem:

$$\begin{aligned} \forall \mathbf{x}, \mathbf{y} \in \text{dom } f, \exists \mathbf{z} \in (\mathbf{x}, \mathbf{y}) \\ f(\mathbf{y}) &= f(\mathbf{x}) + f'(\mathbf{x})(\mathbf{y}-\mathbf{x}) + \frac{1}{2} \underbrace{f''(\mathbf{z})}_{\geq 0} \underbrace{(\mathbf{y}-\mathbf{x})^2}_{\geq 0} \\ \Rightarrow f(\mathbf{y}) &\geq f(\mathbf{x}) + f'(\mathbf{x}) \cdot (\mathbf{y}-\mathbf{x}) \end{aligned}$$

#

Proof sketch for high-dimensional cases:

We consider any 1D line in $\text{dom } f$. If a function cut by any line in the domain is convex, then the function is convex.

Convexity of MSE

First-order derivative

$$J(\mathbf{w}) = \frac{1}{2} \sum_{m=1}^M \left(\sum_{k=0}^d w_k \cdot x_k^{(m)} - t^{(m)} \right)^2$$

$$\begin{aligned} \frac{\partial}{\partial w_i} J(\mathbf{w}) &= \frac{1}{2} \sum_{m=1}^M \cdot 2 \left(\sum_{k=0}^d w_k \cdot x_k^{(m)} - t^{(m)} \right) \cdot x_i^{(m)} \\ &= \sum_{m=1}^M \left(\sum_{k=0}^d w_k \cdot x_k^{(m)} - t^{(m)} \right) \cdot x_i^{(m)} \end{aligned}$$

From the element-wise derivative, we have

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \underbrace{\mathbf{X}^T}_{(d+1) \times M} \underbrace{(\mathbf{X}\mathbf{w} - \mathbf{t})}_{M \times 1} \in \mathbb{R}^{d+1}$$

Noticing that $\begin{pmatrix} | & | \\ x^{(1)} & \cdots & x^{(M)} \\ | & \cdots & | \\ \mathbf{X}^T & & \mathbf{X}\mathbf{w} - \mathbf{t} \end{pmatrix}$ has the

implicit summation over M .

The gradient can also be obtained from

$$J(w) = \frac{1}{2} \|Xw - t\|^2$$

by matrix calculus.

$$\begin{aligned}\nabla_w J(w) &= \nabla_w \left[\frac{1}{2} (Xw - t)^T (Xw - t) \right] \\ &= \nabla_w \left[\frac{1}{2} w^T X^T X w - t^T X w + \frac{1}{2} t^T t \right] \\ &= X^T X w - X^T t\end{aligned}$$

Second-order derivative

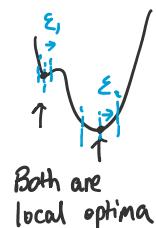
$$\nabla^2 J(w) = X^T X \succeq 0 \quad (\text{PSD})$$

Conclusion: Linear regression is convex in the parameters

Optimality of a convex function

Local optimum

x is a local optimum of $f \iff \exists \varepsilon > 0, \forall z \in \text{dom } f, \text{ if } \|z - x\| < \varepsilon \text{ then } f(x) \leq f(z)$



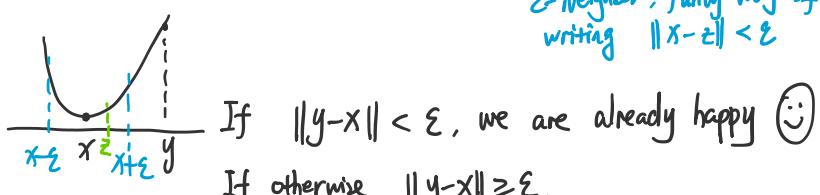
Global optimum

x is a global optimum of $f \iff \forall z \in \text{dom } f, f(x) \leq f(z)$

Theorem: Convex \Rightarrow local optimum is global

Proof. We pick such ε that $\forall z \in N_\varepsilon(x), f(x) \leq f(z)$

ε -neighbor, fancy way of writing $\|x - z\| < \varepsilon$



$$\text{we define } \lambda = \frac{\varepsilon}{2\|y - x\|} \in (0, 1)$$

$$\text{Define } z = (1-\lambda)x + \lambda y$$

$$\|z-x\| = \lambda \|x-y\| = \frac{\varepsilon}{2} < \varepsilon \Rightarrow f(x) \leq f(z)$$

By convexity

$$(1-\lambda)f(x) + \lambda f(y) \geq f((1-\lambda)x + \lambda y)$$

$$= f(z)$$

$$\geq f(x)$$

Thus $f(y) \geq f(x)$. We are also happy 😊 #

Further reading on convexity analysis in unconstrained scenario:

http://www.princeton.edu/~aaa/Public/Teaching/ORF523/S16/ORF523_S16_Lec7_gh.pdf

Closed-form solution to linear regression

$$\text{Set } \nabla_w J(w) = 0$$

$$X^T X w - X^T t = 0$$

$$w = (X^T X)^{-1} X^T t \text{ is the global optimum}$$

Here, we assume $X^T X$ is invertible. If not use pseudo-inverse.

See A.5.4 in https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf

Complexity: $\mathcal{O}(n^{2.3}) \sim \mathcal{O}(n^3)$

https://en.wikipedia.org/wiki/Computational_complexity_of_mathematical_operations#endnote_blockinversion

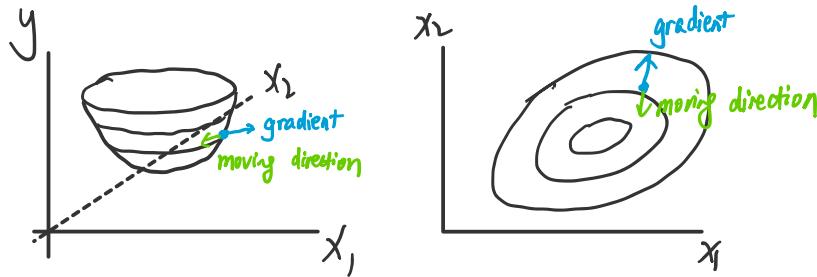
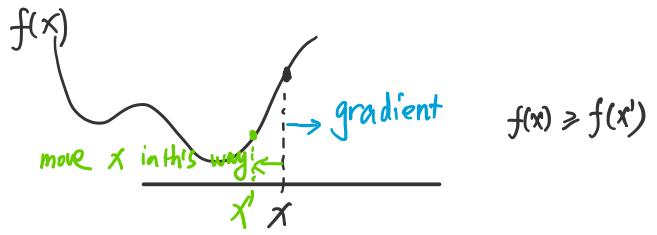
Numerical solution - Gradient descent

If we move a **small** step towards the opposite direction of gradient, Then, the function value decreases.

As $\alpha \rightarrow 0$, let $\mathbf{x}' = \mathbf{x}$

$$\begin{aligned} f(\mathbf{x}') &= f(\mathbf{x}) + [\nabla_{\mathbf{x}} f(\mathbf{x})]^T [-\alpha \nabla_{\mathbf{x}} f(\mathbf{x})] + o(\alpha^2) \\ &= f(\mathbf{x}) - \alpha \|\nabla_{\mathbf{x}} f(\mathbf{x})\|^2 + o(\alpha^2) \end{aligned}$$

$$\Rightarrow f(\mathbf{x}') \leq f(\mathbf{x})$$



Iterative update: We compute the gradient and move a small step against the gradient. Repeat the above criteria for many times.

Gradient Descent Algorithm

Initialize $\mathbf{w} = \mathbf{w}^{(0)}$ randomly

Loop over epochs $t = 0, 1, 2, \dots$:

Compute gradient $\nabla J(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}^{(t)}}$

Update parameters

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} - \alpha \cdot \nabla J(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

Until stopping criterion satisfies

Newton's Method

2nd-order Taylor expansion is more precise than 1st-order.

$$f(\mathbf{x} + \mathbf{t}) = f(\mathbf{x}) + [\nabla f(\mathbf{x})]^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{H}(\mathbf{x}) \mathbf{t} + O(\|\mathbf{t}\|^3)$$

Suppose f is indeed quadratic, then \mathbf{t} can be solved by

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{t}} &= \nabla f(\mathbf{x}) + \mathbf{H}(\mathbf{x}) \mathbf{t} \stackrel{\text{set}}{=} \mathbf{0} \\ \mathbf{t} &= -[\mathbf{H}(\mathbf{x})]^{-1} \nabla f(\mathbf{x}) \end{aligned}$$

However, the objective function J may not be quadratic in \mathbf{w} , and thus multiple iterations are needed.

Newton's Method

Initialize $\mathbf{w} = \mathbf{w}^{(0)}$ randomly

Loop over epochs $t = 0, 1, 2, \dots$:

Compute gradient and Hessian

$$\mathbf{g} = \nabla_{\mathbf{w}} J(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}^{(t)}}, \quad \mathbf{H} = \nabla^2 J(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

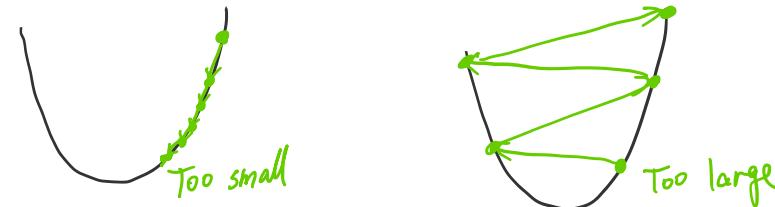
Update parameters $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \mathbf{H}^{-1} \mathbf{g}$

Until stopping criterion satisfies

Comments:

- Initialization of $\mathbf{w}^{(0)}$:
 - Convex: Doesn't matter that much.
 - Non-convex: small random values around the origin $\mathbf{0}$
- Iteration over data
 - One epoch means one iteration over all data samples
 - Stopping criteria
 - Set MaxIter due to budget constraint
 - Until the convergence of J , i.e.,

$$|J(\mathbf{w}^{(t)}) - J(\mathbf{w}^{(t-1)})| < \text{threshold}$$
 - Early stop by validation (TBD in future lectures)
- α is the learning rate
 - Theoretical result: Gradient descent guarantees to converge with some assumptions: <https://www.stat.cmu.edu/~ryantibs/convexopt-F13/scribes/lec6.pdf>
 - In practice:
 - Too small \Rightarrow slow, and performance may not be very good
 - Too large \Rightarrow overshoot



- Could be tuned in a multiplicative way, e.g.,

$$\dots, 0.01, 0.03, 0.1, 0.3, 1, 3, \dots$$
- Learning rate decay
 - Start with a relatively large α

- Decrease α during training
- Recent advances: Adaptive learning rate for different parameters

- Gradient

- Full-batch gradient descent

$$J = \frac{1}{M} \sum_{i=1}^M J^{(i)}$$

- Mini-batch gradient descent

In each iteration in the loop,

$$J = \frac{1}{B} \sum_{b=1}^B J^{(n_b)}$$

where $B < N$, and $n_1, n_2, \dots, n_B \in \{x_1, \dots, x_M\}$,

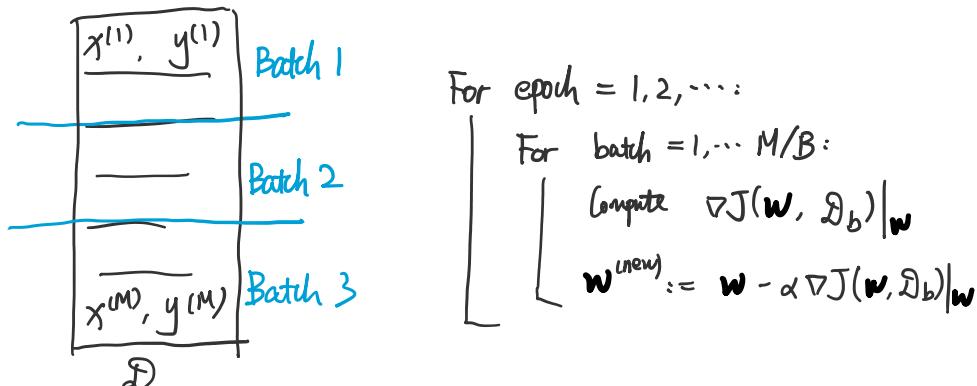
- Stochastic gradient descent (SGD)

Mini-batch with $B = 1$

Convergence: <https://www.cs.ubc.ca/~schmidtm/Courses/540-W19/L11.pdf>

Essentially, a batch is a sampled subset of the entire dataset. Two settings:

- Sample with replacement
- Sample without replacement: In this case, all batches form a partition of the dataset



Pros and Cons

	Full batch	SGD ($B = 1$)
Pros	Exact gradient	Gradient could be noisy
Cons	Slow in each param. update	Efficient for each update

Mini-batch gradient descent is a tradeoff. Besides efficiency concerns, evidence shows that small batch gradient descent is better (yields more generalizable solutions) than full batch. The theoretical reasons are not completely clear.

Comparing closed-form solution and numerical solution

- For linear regression:
 - Closed-form solution converges better (guaranteed global optimum) and is faster if the number of sample is not too large.
 - Numerical solution works regardless of the existence of matrix inverse. It is more efficient if we have too many samples.
- For other optimization/computation:
 - Having a closed-form solution is rare, but should always be considered before trying numerical methods. In gradient descent, for example, we compute the gradient by closed-form formula (although the optimum is not closed-form), instead of numerical gradient computing.
 - Numerical methods are more flexible. Gradient descent works for not only convex problems, but also non-convex ones. In non-convex gradient descent, we may be stuck at a local optimum. Restarting gradient descent by different initial parameters may alleviate this problem.