

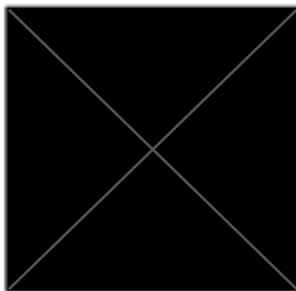
Spring 26 IntroToML Lab 2

Convexity, Gradient Descent, and Logistic Regression

Patrick Shen

NYU Center for Data Science

January 30, 2026



Agenda

- Review of GD and SGD
- GD on logistic regression
- convexity theory \rightarrow Lec 4.
- Notebook demo

Review of Gradient Descent

Consider empirical risk

$$f(x) = \frac{1}{\underbrace{n}} \sum_{i=1}^n f_i(x), \quad x \in \mathbb{R}^d.$$

GD update:

$$\theta_k = \underbrace{(\theta_{k-1})} - \eta_k \underbrace{\nabla_{\theta_{k-1}}}_{\uparrow} \underbrace{f(x)}.$$

If $\theta_k \cong \underbrace{\theta_{k-1}}$, then $\underbrace{\theta_{k-1}}$ is a fixed point of the update, and (for convex f) a global minimizer.

Stochastic Gradient Descent

SGD update:

$$\theta_k = \theta_{k-1} - \eta_k \underbrace{\nabla f}_{\theta_{k-1}}(\underbrace{x}_{i_k}),$$

where i_k is sampled uniformly from $\{1, 2, \dots, n\}$.

Unbiasedness:

$$\frac{1}{n}$$

$$\mathbb{E}[\underbrace{\nabla f}_{\theta_{k-1}}(\underbrace{x}_{i_k})] = \underbrace{\nabla_{\theta_{k-1}} f(x)}.$$

\Rightarrow mini-batch. $|B| = [1, \dots, n]$

Logistic Regression: model

softmax

Outputs are in $[0, 1]$.

Let $\underline{X} \in \mathbb{R}^{m \times d}$, $\underline{y} \in \{0, 1\}^m$, parameters $\underline{w} \in \mathbb{R}^d$, $\underline{b} \in \mathbb{R}$.



$$\text{n. of sample } \underline{z} = \underline{X}\underline{w} + \underline{b}\mathbf{1}, \quad \hat{y} = \sigma(z), \quad \sigma(z) = \frac{1}{1 + e^{-z}}.$$

$\in \mathbb{R}^{m \times 1}$

Logistic loss

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] .$$

Logistic loss

$$\underline{J(w, b)} = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)].$$

Exercise.

Hint: Chain rule

- Write $\underline{\frac{\partial J}{\partial z}}, \underline{\frac{\partial J}{\partial w}}, \underline{\frac{\partial J}{\partial b}}$.
- Write one step of GD for (w, b) .

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell_i$$

$$\ell_i = -[y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)] \quad \hat{y}_i = \sigma(z_i) = \frac{1}{1 + e^{-z_i}}$$

$$\frac{\partial \ell_i}{\partial z_i} = \frac{\partial \ell_i}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i}$$

$$\textcircled{1} \frac{\partial \ell_i}{\partial \hat{y}_i} = -\left(\frac{y_i}{\hat{y}_i} - \frac{1-y_i}{1-\hat{y}_i}\right)$$

$$\textcircled{2} \frac{\partial \hat{y}_i}{\partial z_i} = \frac{0(1+e^{-z_i}) - (1+e^{-z_i}) \cdot 1}{(1+e^{-z_i})^2}$$

$$= \frac{-(-1)e^{-z_i}}{(1+e^{-z_i})} \cdot \frac{1}{1+e^{-z_i}}$$

$$= (1-\hat{y}_i)\hat{y}_i$$

$$\textcircled{1} \times \textcircled{2} = -\left(\frac{y_i(1-\hat{y}_i) - \hat{y}_i(1-y_i)}{\hat{y}_i(1-\hat{y}_i)}\right) \cdot (1-\hat{y}_i)\hat{y}_i = \hat{y}_i(1-y_i) - y_i(1-\hat{y}_i)$$

$$= \hat{y}_i - \hat{y}_i y_i - y_i + y_i \hat{y}_i$$

$$= \hat{y}_i - y_i$$

$$\frac{\partial J(w, b)}{\partial z} = \left(\frac{1}{m}\right)(\hat{y} - y), \quad \hat{y}, y \in \mathbb{R}^m$$

$\rightarrow \mathbb{R}^{m \times 1}$ $\rightarrow x_i \in \mathbb{R}^{d \times 1}$

$$\text{II: } z = Xw + b \quad , \quad z_i = x_i^T w + b$$

$$\frac{\partial z_i}{\partial w} = x_i \quad \frac{\partial z_i}{\partial b} = 1$$

$$\frac{\partial J(w, b)}{\partial w} = \frac{\partial J(w, b)}{\partial z} \cdot \frac{\partial z}{\partial w} = \left(\frac{1}{m}\right) X^T (\hat{y} - y) \quad (\hat{y} - y) \in \mathbb{R}^{m \times 1}$$

$X^T \mathbb{R}^{d \times m} \cdot \mathbb{R}^{m \times 1} \rightarrow \mathbb{R}^{d \times 1}$

$$\frac{\partial J(w, b)}{\partial b} = \frac{\partial J(w, b)}{\partial z} \cdot \frac{\partial z}{\partial b} = \frac{1}{m} \mathbf{1}^T (\hat{y} - y)$$

$$w_k = w_{k-1} - \eta \frac{\partial J(w, b)}{\partial w_{k-1}}$$

$$b_k = b_{k-1} - \eta \frac{\partial J(w, b)}{\partial b_{k-1}}$$

Gradients and one-step GD (standard form)

For logistic regression with $\hat{y} = \sigma(z)$ and $z = Xw + b\mathbf{1}$:

$$\frac{\partial J}{\partial z} = \frac{1}{m}(\hat{y} - y), \quad \frac{\partial J}{\partial w} = \frac{1}{m}X^\top(\hat{y} - y), \quad \frac{\partial J}{\partial b} = \frac{1}{m}\mathbf{1}^\top(\hat{y} - y).$$

Gradients and one-step GD (standard form)

For logistic regression with $\hat{y} = \sigma(z)$ and $z = Xw + b\mathbf{1}$:

$$\frac{\partial J}{\partial z} = \frac{1}{m}(\hat{y} - y), \quad \frac{\partial J}{\partial w} = \frac{1}{m}X^\top(\hat{y} - y), \quad \frac{\partial J}{\partial b} = \frac{1}{m}\mathbf{1}^\top(\hat{y} - y).$$

One GD step:

$$w_k = w_{k-1} - \eta_k \frac{\partial J}{\partial w_{k-1}}, \quad b_k = b_{k-1} - \eta_k \frac{\partial J}{\partial b_{k-1}}.$$

Notebook exercise

Optimization: the key question

The key question in optimization is to find the smallest (or largest) value of a function.

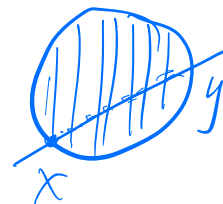
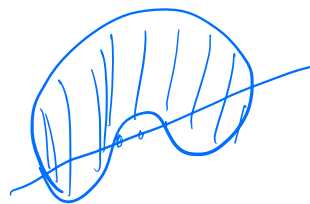
$$\min f(x) \quad \text{s.t.} \quad x \in \Omega \subseteq \mathbb{R}^d,$$

where:

- Ω is the **domain / feasible set**
- $f(x)$ is the **objective function**

Definition (Convex set). A set Ω is convex if for all $x, y \in \Omega$ and all $\lambda \in [0, 1]$,

$$\lambda x + (1 - \lambda)y \in \Omega.$$



Definition (Convex set). A set Ω is convex if for all $x, y \in \Omega$ and all $\lambda \in [0, 1]$,

$$\lambda x + (1 - \lambda)y \in \Omega.$$

Exercise. Let

$$C = \{x \in \mathbb{R}^n : \|x\|_2 \leq 1\}.$$

Prove that C is a convex set.

triangle inequality $\| \lambda x + (1 - \lambda)y \|_2 \leq \lambda \|x\|_2 + (1 - \lambda) \|y\|_2 = 1$

Characterization of convexity

Theorem. The following statements are equivalent:

- 1 f is convex.
- 2 For any $\lambda \in [0, 1]$ and any x, y in the domain of f ,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

First-order characterization of convexity

Assume $\underline{f} : \Omega \rightarrow \mathbb{R}$ is **differentiable** and the domain $\underline{\Omega}$ is convex.

Theorem (First-order condition). The following are equivalent:

- ① f is convex.
- ② For all $x, y \in \Omega$,

$$\underline{f(y)} \geq \underline{f(x)} + \langle \nabla \underline{f(x)}, y - x \rangle.$$

First-order characterization of convexity

Assume $f : \Omega \rightarrow \mathbb{R}$ is **differentiable** and the domain Ω is convex.

Theorem (First-order condition). The following are equivalent:

- ① f is convex.
- ② For all $x, y \in \Omega$,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle.$$



Geometric meaning: the tangent hyperplane at x is a global underestimator of f .

Second-order characterization of convexity

Assume $f : \Omega \rightarrow \mathbb{R}$ is **twice differentiable** and the domain Ω is convex.

Theorem (Second-order condition). The following are equivalent:

- 1 f is convex.
- 2 For all $x \in \Omega$, the Hessian is positive semidefinite:

$$\nabla^2 f(x) \succeq 0 \iff v^\top \nabla^2 f(x) \underline{v} \geq 0 \quad \forall v \in \mathbb{R}^d.$$

Special case: if $\nabla^2 f(x) \succ 0$ for all x , then f is **strictly convex**.

Useful lemmas: why GD can work (intuition)

Lemma 1. For a convex function $f(x)$, all local minima (if they exist) are global minima.

- Proof idea: contradiction.

Lemma 2. Let $\{g_i\}_{i=1}^n$ be convex functions. Then

$$g(x) := \frac{1}{n} \sum_{i=1}^n \underbrace{g_i(x)}$$

is convex.

First-order optimality for differentiable convex f

Theorem. If f is differentiable and convex, then the following are equivalent:

- ① x^* is a (local/global) minimizer of f
- ② $\nabla f(x^*) = 0$

Exercise: quadratic objective

$$\frac{1}{2} w^T w = \frac{1}{2} \sum_{i=1}^m w_i^2 \stackrel{d}{=} \frac{1}{2} \sum_{i=1}^m 2w_i = \sum_{i=1}^m w_i = W^T \cdot \mathbf{1}.$$

Let

$$f(w) = \frac{1}{2} \|w\|_2^2, \quad w \in \mathbb{R}^d.$$

- 1 Compute $\nabla f(w)$.
 - 2 Write the gradient update $w_k = w_{k-1} - \eta \nabla f(w_k)$.
- $$\frac{d}{dw} \frac{1}{2} w^T w = \frac{1}{2} \times 2 \times w = w.$$

Exercise: convexity of logistic loss

Exercise. Show that the logistic loss function $J(w, b)$ is convex.

Hint. Compute the Hessian and show it is positive semidefinite (PSD).

Wrap-up

- Convex sets/functions and key properties ← Lecture 4
- GD and SGD updates for empirical risk minimization
- Logistic regression: model, loss, gradients, and GD steps

