

Admin:

→ Quiz 2 grades are out *You did well!* 😊

Today: Applications of gradient descent

- Recap of gradient descent method
- Linear regression (with mean squared error loss)
- Solving regression problems via gradient descent
- Linear system solving
- Nonlinear regression via monomial dimension expansion

Last time:Problem: Unconstrained Minimization

Given $f: \mathbb{R}^n \rightarrow \mathbb{R}$ f is continuous & smooth

find

$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$$

(provided minimum of f exists)Algorithm: Gradient descent (GD) method
(locally greedy approach)→ Start with some $x^{(0)}$ → For $i = 0, 1, \dots, T$

$$x^{(i+1)} = x^{(i)} - \underbrace{(\eta_i)}_{\text{Step size}} \underbrace{\nabla f(x^{(i)})}_{\text{Gradient}}$$

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

(2)

Intuition:

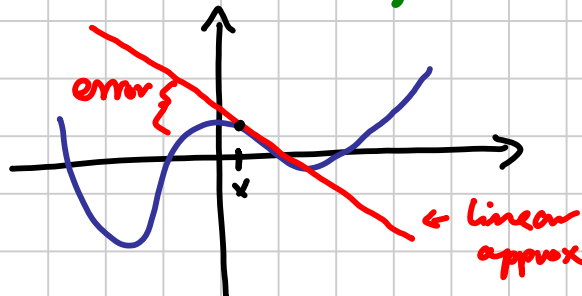
$$f(x+\delta) = \underbrace{f(x) + [\nabla f(x)]^T \delta}_{\text{linear approx.}} + \underbrace{\frac{1}{2} \delta^T \nabla^2 f(x) \delta}_{\text{"error" term}} + \dots$$

Hessian

(GD minimizes that)

"error" term
(we choose δ_i
to be small enough
so our progress
is not cancelled
by these terms)

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_n \partial x_n} \end{bmatrix} = \left[\frac{\partial^2 f(x)}{\partial x_i \partial x_j} \right]_{ij}$$



Convergence of GD:

For $T \rightarrow \infty$, $x^{(T)}$ either diverges to $-\infty$, i.e.,

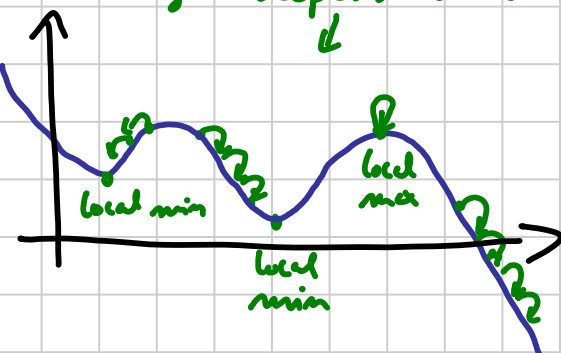
$$\lim_{T \rightarrow \infty} f(x^{(T)}) = -\infty$$

where we end up
might depend on starting point!

Or converges to a critical point, i.e.,

$$\lim_{T \rightarrow \infty} x^{(T)} = \hat{x} \text{ s.t. } \nabla f(\hat{x}) = 0$$

$\Rightarrow \hat{x}$ is a local minimum,
maximum, or a saddle point

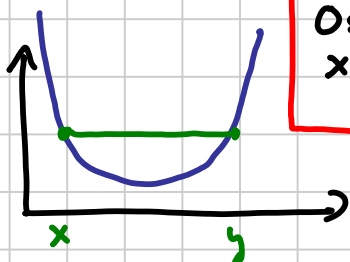


divergence

Convexity: Function f is convex iff

$$\forall 0 \leq \lambda \leq 1 \\ \forall x, y \in \mathbb{R}^n$$

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



Equivalently,

$$\forall x, \delta \in \mathbb{R}^n \quad f(x+\delta) \geq f(x) + \nabla f(x)^T \delta$$

Key observation: If f is convex then (3)
every critical point is a global minimum, i.e.,
 $\nabla f(\tilde{x}) = 0 \Leftrightarrow \tilde{x}$ is a global minimum.
← this is the key direction

So, if f is convex & global minimum exists

\Rightarrow GD provably finds it

(Alas, many important problems correspond to f that are NOT convex.)

Convergence analysis: (for (strongly) convex functions)

Def. f is β -smooth iff

$$\forall x, \delta \in \mathbb{R}^n$$

$$\delta^T \nabla^2 f(x) \delta \leq \beta \|\delta\|^2$$

for some $\beta > 0$

Def. f is α -strongly convex iff

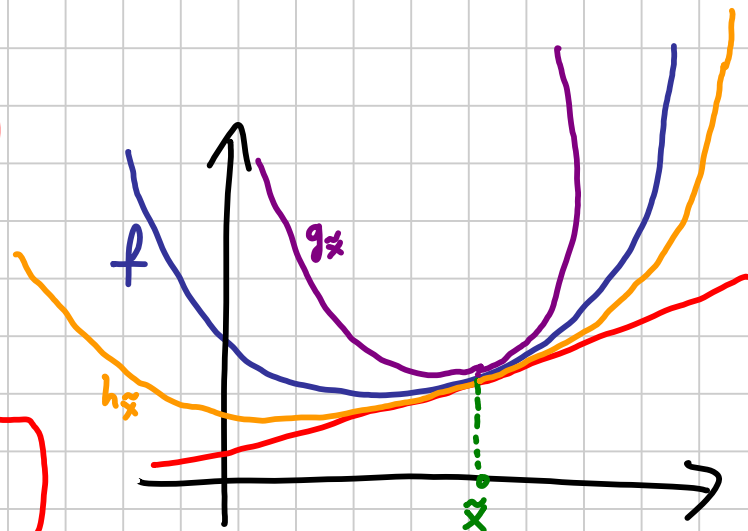
$$\forall x, \delta \in \mathbb{R}^n$$

$$\delta^T \nabla^2 f(x) \delta \geq \alpha \|\delta\|^2$$

for some $\alpha > 0$

("normal" convexity corresponds to $\alpha = 0$)

\Rightarrow Always $\alpha \leq \beta$



$$g_{\tilde{x}}(x) = f(\tilde{x}) + \nabla f(\tilde{x})^T (x - \tilde{x}) + \frac{\beta}{2} \|x - \tilde{x}\|^2$$

(upper bounding " β -smoothness" parabola)

$$\Rightarrow f(x) \leq g_{\tilde{x}}(x) \quad \forall x, \tilde{x}$$

(lower bounding " α -strong conv." parabola) $\Rightarrow f(x) \geq h_{\tilde{x}}(x) \quad \forall x, \tilde{x}$

$$h_{\tilde{x}}(x) = f(\tilde{x}) + \nabla f(\tilde{x})^T (x - \tilde{x}) + \frac{\alpha}{2} \|x - \tilde{x}\|^2$$

Theorem (proved last time)

(4)

If f is β -smooth & α -strongly convex then for any $\varepsilon > 0$

$$f(x^{(T)}) - f(x^*) \leq \varepsilon \text{ whenever } T = \Omega\left(K \log \frac{f(x^{(0)}) - f(x^*)}{\varepsilon}\right)$$

where $K = \frac{\beta}{\alpha}$ is the condition number of f

initial suboptimality

In other words: For such f , GD gives an ε -approx solution after $\leq O\left(K \log \frac{f(x^{(0)}) - f(x^*)}{\varepsilon}\right)$ iterations

Intuition on the dependence on K :

$\rightarrow K$ measures the quality of the local approximation of f at \bar{x} via upper bounding "parabola" $g_{\bar{x}}$ and lower bounding parabola $h_{\bar{x}}$

Geometric intuition:

Consider $f(x) = \frac{a}{2} x_1^2 + \frac{1}{2} x_2^2$ for $a \geq 1$

$$\nabla f(x) = \begin{bmatrix} a \\ 1 \end{bmatrix}$$

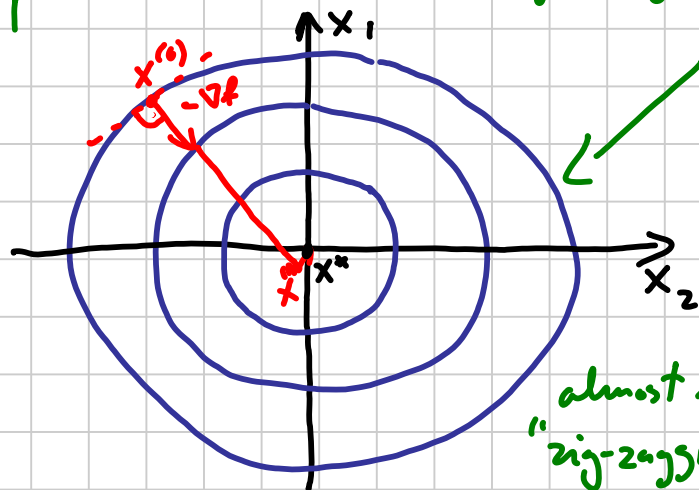
$$\nabla^2 f(x) = \begin{bmatrix} a & 0 \\ 0 & 1 \end{bmatrix}$$

$\Rightarrow f$ is a -smooth & 1-strongly convex

$\Rightarrow K = a \quad \Rightarrow x^* = (0, 0)$

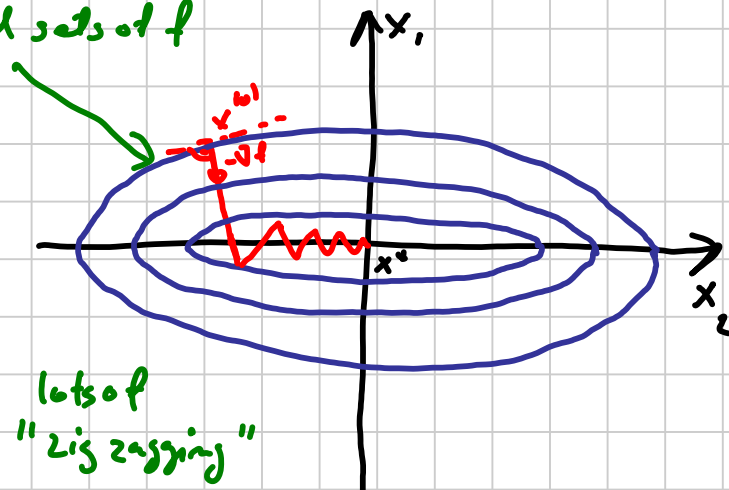
Visualization of a trajectory of GD starting from some $x^{(0)}$

$K \approx 1$



almost no
"zig-zagging"

level sets of f



lots of
"zig-zagging"

Key application domain of gradient descent is training machine learning models

(5)

Illustrative example: Linear regression

Input: m datapoints $x^{(1)}, \dots, x^{(m)} \in \mathbb{R}^n$ & m labels $y^{(1)}, \dots, y^{(m)} \in \mathbb{R}$

Goal: Find a linear function h_w that "predicts" each $y^{(j)}$ given $x^{(j)}$, i.e., we want

$$h_w(x^{(j)}) = \sum_{i=1}^n w_i \cdot x_i^{(j)} \approx y^{(j)} \quad \text{for each } j=1, \dots, m$$

Think: Predict the score of 6.046 students on the final exam based on their performance during the semester

→ each data point $x^{(j)}$ corresponds to 6.046 student

→ each coordinate of $x^{(j)}$ encodes score on one of the psets or the quizzes

→ each $y^{(j)}$ is the score on the final

Note: Ulog, we allow our function h_w to have a constant term → just make one "dummy" coordinate of datapoints $x^{(1)}, \dots, x^{(m)}$ be always equal to 1

So, our goal is to compute $w \in \mathbb{R}^n$ s.t. the resulting linear function

$$h_w(x) = \sum_{i=1}^n w_i x_i = w^T x$$

has "best fit" to data points

⇒ good fit indicates underlying linear relationship

How to quantify that?

⇒ hope that such h_w "generalizes" to yet unseen data

Popular choice: Mean squared error (MSE)

6

$$L(u) = \frac{1}{n} \sum_{j=1}^m E_j(u)^2$$

where

$$E_j(u) = h_u(x^{(j)}) - y^{(j)}$$

= "residual error" on $x^{(j)}$, i.e. the difference between h_u "prediction" for $x^{(j)}$ and the actual value of $y^{(j)}$

Our goal:

Compute

$$u^* = \underset{u \in \mathbb{R}^n}{\operatorname{argmin}} L(u)$$

(Corresponds to a max likelihood estimate if we assume that the data set has linear dependencies + Gaussian noise)

\Rightarrow This is an unconstrained minimization problem!

\Rightarrow We can just use GD

Observe:

$$\nabla E_j(u) = \begin{bmatrix} x_1^{(j)} \\ \vdots \\ x_n^{(j)} \end{bmatrix}$$

Remember: Now, u are the variables (not x)

$$\nabla L(u) = \frac{2}{n} \sum_{j=1}^m E_j(u) \cdot \nabla E_j(u) = \frac{2}{n} \sum_{j=1}^m E_j(u) \cdot x^{(j)}$$

\Rightarrow Each datapoint $x^{(j)}$ contributes $\frac{2}{n} E_j(u) x^{(j)} = v_j$ to the gradient

\Rightarrow update $u' \leftarrow u - \eta v_j$

reduces the error on data point $x^{(j)}$ since

$$E_j(u') = (u')^T x^{(j)} - y^{(j)} = (u - \eta v_j)^T x^{(j)} - y^{(j)} =$$

$$= E_j(u) - \eta v_j^T x^{(j)} = E_j(u) - \frac{2}{n} \eta E_j(u) \|x^{(j)}\|^2 = E_j(u) \left(1 - \frac{2}{n} \eta\right)$$

But could increase error on other example

$\frac{1}{n} L^1 \rightarrow \frac{1}{n} L^2$
 \uparrow
 $\frac{1}{n} L^1 \rightarrow \frac{1}{n} L^2$

⇒ Gradient update: $-\eta \nabla L(u) = -\eta \sum v_j$ ⑦
 takes a "majority" vote of all data point updates

How fast will it converge (if at all)?

Observe:

$$\nabla^2 L(u) = \frac{1}{m} \sum_{j=1}^m \nabla^2 (E_j(u)^2) = \frac{1}{m} \sum_{j=1}^m \left[\frac{\partial^2 (E_j(u)^2)}{\partial k \partial l} \right]_{k,l}$$

$$\Rightarrow \nabla^2 L(u) = \frac{2}{m} \sum_{j=1}^m \underbrace{x^{(j)} (x^{(j)})^T}_{\substack{\uparrow \\ \text{this is rank-one} \\ n \times n \text{ matrix}}} \quad \quad \quad 2 \frac{\partial^2 (E_j(u) x_L^{(j)})}{\partial k \partial l}$$

$$\quad \quad \quad 2 x_k^{(j)} x_l^{(j)}$$

Recall: For any vector $v \in \mathbb{R}^n$
 the (rank-one) matrix vv^T is positive semi-definite
 i.e. $\delta^T (vv^T) \delta \geq 0$ for any $\delta \in \mathbb{R}^n$ (PSD)

⇒ $\nabla^2 L(u)$ is PD too

Means that $\nabla^2 L(u)$
 is positive definite
 (PD)

⇒ $L(u)$ is convex

In fact: as long as $m \geq n$ & all $x^{(j)}$'s are linearly indep.

$$\delta^T \nabla^2 L(u) \delta \geq \alpha \|\delta\|^2 \text{ for some } \alpha > 0$$

⇒ $L(u)$ is α -strongly convex

Also: We always have some β s.t.

$$\delta^T \nabla^2 L(u) \delta \leq \beta \|\delta\|^2 \text{ for some } \beta > 0$$

⇒ $L(u)$ is β -smooth

⇒ u^* exists & is unique ⇒ GD gets ϵ -close to it after $O\left(\frac{\beta}{\alpha} \log \frac{L(u) - L(u^*)}{\epsilon}\right)$ iter

In fact: We can get an explicit formula for u^*
(this is quite unusual)

8

Recall: u^* has to be a critical point, i.e., we must have

$$\nabla L(u^*) = 0$$

← Normal equations for L

$$\Rightarrow \sum_{j=1}^m E_j(u^*) x^{(j)} = 0 \quad | \cdot \frac{1}{2}$$

$$\Rightarrow \sum_{j=1}^m ((u^*)^T x^{(j)} - y^{(j)}) x^{(j)} = 0$$

$$\Rightarrow \sum_{j=1}^m (u^*)^T x^{(j)} \cdot x_i^{(j)} = \sum_{j=1}^m y^{(j)} x_i^{(j)} \quad \forall i=1, \dots, n$$

$$\Rightarrow (X^T X) u^* = X^T y$$

where

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}} \right\}^n$$

$$\Rightarrow u^* = (X^T X)^{-1} X^T y$$

Data matrix

Note:
$$\begin{aligned} \nabla L(u) &= \frac{2}{m} \sum_{j=1}^m x^{(j)} x^{(j)T} = \\ &= \frac{2}{m} X^T X \end{aligned}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}} \right\}^m \quad \text{Label vector}$$

\Rightarrow $X^T X$ invertible (and u^* well-defined)
iff $\nabla^2 L(u)$ is PD iff $L(u)$ α -strongly convex for $\alpha > 0$

↳ This is a "sanity check", as u^* has to be defined uniquely if $L(u)$ α -strongly convex

So, we can use GD to solve (linear) regression but what else?

⑨

Solving linear systems:

Given an invertible $n \times n$ matrix A & vector $b \in \mathbb{R}^n$,

find $x^* \in \mathbb{R}^n$ s.t.

$$Ax = b$$

→ Classic approach:

Compute $x^* = A^{-1}b$ directly via, e.g., Gaussian elimination

Problems: → Always fairly slow $O(n^3)$ ≈ 2.373
(even if matrix is sparse or "nice")

→ Numerically problematic
(dividing by small numbers is bad)

→ Iterative approach:

→ Start with some $x^{(0)}$

→ Iteratively improve your solution $x^{(i)}$ to get solution $x^{(i+1)}$

How to do this? With gradient descent!

Consider a function:

$$f_A(x) = \frac{1}{2} x^T (A^T A) x - A^T b$$

Observe: $\nabla f_A(x) = (A^T A)x - A^T b$

← cheap to compute:
 $O(n + \# \text{ of non-zero entries of } A)$
small if A is sparse

Algorithm:

→ Start with some $x^{(0)}$, e.g., $x^{(0)} = 0$

→ For $j = 0, 1, \dots, T$

$$x^{(j+1)} = x^{(j)} - \eta_j \nabla f_A(x^{(j)}) = x^{(j)} - \eta_j A^T \underbrace{(Ax^{(j)} - b)}_{\text{residual error of } x^{(j)}}$$

Why does it work?

Note: For any critical point \hat{x} we have that

$$\nabla f_A(\hat{x}) = 0 \Leftrightarrow A^T(A\hat{x} - b) = 0$$

$$\Leftrightarrow \underbrace{A\hat{x} = b}$$

(we use the fact that A is invertible)

\Rightarrow Any critical point is a solution

In fact:

$$\nabla^2 f_A(x) = A^T A$$

$\Rightarrow \nabla^2 f_A(x)$ is PD if A is invertible (for any x)

$\Rightarrow f_A$ is α -strongly convex for some $\alpha > 0$ (& also β -smooth)

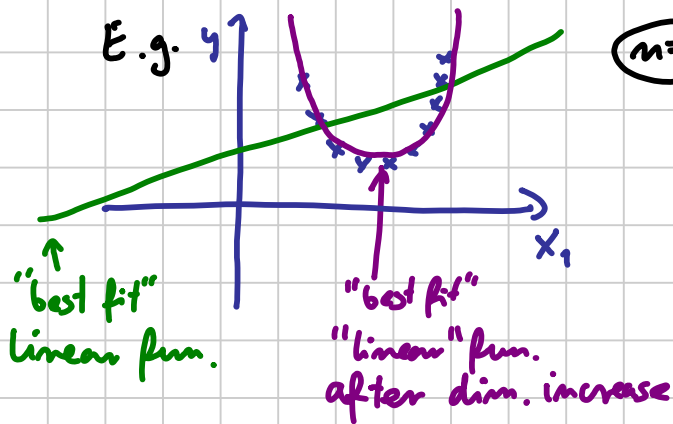
\Rightarrow Solution x^* is unique and GD will compute an ε -approx. solution after $\underbrace{O\left(\frac{\beta}{\alpha} \log \frac{f_A(0) - f_A(x^*)}{\varepsilon}\right)}_{\text{iterations}}$

Advantage: Condition number $\kappa = \frac{\beta}{\alpha}$ can often be (made) small

Beyond Convexity

(11)

So far, our classifier h_u was extremely simple, i.e., linear
Turns out that linear classifiers can still be really powerful!
One just needs to increase the number of dimensions



data points $x^{(j)} \in \mathbb{R}$ & $y^{(j)} = (x^{(j)})^2$
 \Rightarrow impossible to find a line with a good fit

But: If we expand the dimensions via mapping $\varphi: \mathbb{R} \rightarrow \mathbb{R}^2$

$$\varphi(x) = (x, x^2)$$

\Rightarrow "lines" in transformed space correspond to parabolas in the original space

If $z^{(j)} = \varphi(x^{(j)}) \in \mathbb{R}^2$
then a "line" classifier
 $h_u(z) = w_1 \cdot z_1 + w_2 \cdot z_2$

with $w = (0, 1)$ is such that

$$\forall_j \quad h_u(z^{(j)}) = z_2^{(j)} = (x^{(j)})^2 = y^{(j)}$$

← perfect fit! ✓

So, such expanding of the dimension into the $O(m^d)$ -dimensional space of monomials x_i^k for $1 \leq i \leq m$ & $0 \leq k \leq d$ makes "linear" functions express degree-d polynomials

\Rightarrow Significant boost in expressability
(but at the price of blowing up the dimension!)