

CS-GY 6763: Lecture 8

Second Order Conditions, Online and Stochastic Gradient Descent

NYU Tandon School of Engineering, Prof. Christopher Musco

Given a function f to minimize, assume we have:

- **Function oracle:** Evaluate $f(\mathbf{x})$ for any \mathbf{x} .
- **Gradient oracle:** Evaluate $\nabla f(\mathbf{x})$ for any \mathbf{x} .

Goal: Minimize the number of oracle calls to find $\tilde{\mathbf{x}}$ such that $f(\tilde{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon$.

Prototype gradient descent method:

- Choose starting point $\mathbf{x}^{(0)}$.
 - For $i = 0, \dots, T$:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
 - Return $\mathbf{x}^{(T)}$ (or similar).
- step size
learning rate.*

Intuition: Last time we showed that, for sufficiently small η , $f(\mathbf{x}^{(i+1)}) \leq f(\mathbf{x}^{(i)})$. So the algorithm eventually finds a (local) minimum. The question is, how fast.

BASIC GRADIENT DESCENT ANALYSIS

Assume:

- f is convex.
- Lipschitz function: for all \mathbf{x} , $\|\nabla f(\mathbf{x})\|_2 \leq G$.
- Starting radius: $\|\mathbf{x}^* - \mathbf{x}^{(0)}\|_2 \leq R$.

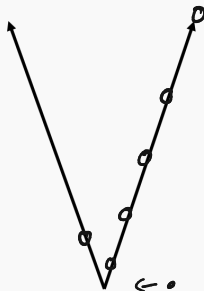
Gradient descent:

- Choose number of steps T .
- Starting point $\mathbf{x}^{(0)}$. E.g. $\mathbf{x}^{(0)} = \vec{0}$.
- $\eta = \frac{R}{G\sqrt{T}}$
- For $i = 0, \dots, T$:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)})$.

BASIC GRADIENT DESCENT ANALYSIS

Claim (GD Convergence Bound)

If we run GD for $T \geq \frac{R^2 G^2}{\epsilon^2}$ iterations then $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$.



Proof was made tricky by the fact that $f(\mathbf{x}^{(i)})$ does not improve monotonically. We can “overshoot” the minimum.

PROJECTED GRADIENT DESCENT

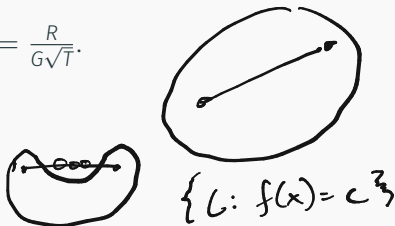
Given function $f(x)$, set \mathcal{S} , and access to **projection oracle**

$$P_{\mathcal{S}}(\underline{x}) = \arg \min_{y \in \mathcal{S}} \|x - y\|_2.$$

$$\min_{x \in \mathcal{S}} f(x)$$

Projected gradient descent:

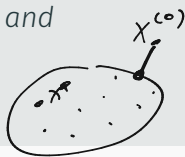
- Select starting point $x^{(0)}$, $\eta = \frac{R}{G\sqrt{T}}$.
- For $i = 0, \dots, T$:
 - $z = x^{(i)} - \eta \nabla f(x^{(i)})$
 - $\underline{x}^{(i+1)} = P_{\mathcal{S}}(z)$
- Return $\hat{x} = \arg \min_i f(x^{(i)})$.



Claim (PGD Convergence Bound)

If f, \mathcal{S} are convex, $\|\nabla f(x)\|_2 \leq \underline{G}$ for all $x \in \mathcal{S}$ and $\|x^{(0)} - x^*\|_2 \leq \underline{R}$.

If $\underline{T} \geq \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{x}) \leq f(x^*) + \epsilon$.



BEYOND THE BASIC BOUND

$$\epsilon = .1 \quad \frac{1}{(10000)^2}$$

The previous bounds are optimal for convex first order optimization in general.

But in practice, the dependence on $1/\epsilon^2$ is pessimistic: gradient descent typically requires far fewer steps to reach ϵ error.

Previous bounds only make a very weak first order assumption:

$$\|\nabla f(x)\|_2 \leq G.$$

In practice, many function satisfy stronger assumptions.

$$|f'(x)| \leq G$$

SECOND ORDER CONDITIONS

Today we will talk about assumptions that involve the second derivative of f .

In particular, we say that a scalar function f is α -strongly convex and β -smooth if for all x :



$$\alpha \leq f''(x) \leq \beta.$$

$$0 \leq f''(x)$$

$$\rightarrow \alpha > 0$$

We will give appropriate generalizations of these conditions to multi-dimensional functions shortly.

Take away: Having either an upper and lower bound on the second derivative helps convergence. Having both helps a lot.

$\frac{1}{\epsilon}$ instead of $1/\epsilon^2$

Take away: Having either an upper and lower bound on the second derivative helps convergence. Having both helps a lot.

Number of iterations for ϵ error:

| | G-Lipschitz \rightarrow β -smooth | |
|-------------------------|---|---|
| R bounded start | $O\left(\frac{G^2 R^2}{\epsilon^2}\right)$ | $\rightarrow O\left(\frac{\beta R^2}{\epsilon}\right)$ |
| α -strong convex | $O\left(\frac{G^2}{\alpha \epsilon}\right)$ | $O\left(\frac{\beta}{\alpha} \underline{\log(1/\epsilon)}\right)$ |

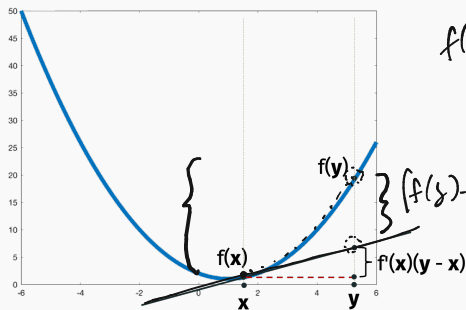
As we defined them so far, smoothness and strong convexity require f to be twice differentiable. On the other hand, gradient descent only requires first order differentiability.

SECOND ORDER CONDITIONS

Equivalent conditions: \rightarrow for y, x

$$\underline{f''(x) \leq \beta} \iff ([f(y) - f(x)] - f'(x)(y-x) \leq \frac{\beta}{2}(y-x)^2)$$

$$f''(x) \geq \alpha \iff [f(y) - f(x)] - f'(x)(y-x) \geq \underline{\underline{\frac{\alpha}{2}(y-x)^2}}$$



$f(y) - f(x) \geq f'(x)(y-x)$
for a convex function.

Recall: For all convex functions $\underline{[f(y) - f(x)] - f'(x)(y-x) \geq 0}$.

SECOND ORDER CONDITIONS

Proof that $f''(x) \leq \beta \Rightarrow [f(y) - f(x)] - f'(x)(y-x) \leq \frac{\beta}{2}(y-x)^2$:

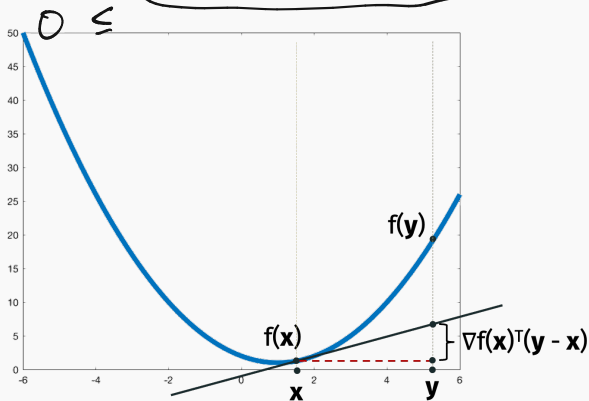
$$\begin{aligned} \underline{f(y) - f(x)} &= \int_x^y f'(t) dt \\ &\leq \int_x^y \underline{f'(x) + (t-x)\beta} dt \\ &= f'(x)(y-x) + \int_x^y \beta(t-x) dt \\ &= \underline{f'(x)(y-x)} + \beta \underbrace{\left. \frac{(t-x)^2}{2} \right|_x^y} \\ &= \underline{\frac{\beta}{2}(y-x)^2} \end{aligned}$$

Proof for α -strongly convex is similar, as are the other directions.

MULTIDIMENSIONAL GENERALIZATION

A function is α -strongly convex and β -smooth if for all \mathbf{x}, \mathbf{y} :

$$\frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \leq [f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$



Definition (β -smoothness)

A function f is β smooth if, for all \mathbf{x}, \mathbf{y}

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$$

I.e., the gradient function is a β -Lipschitz function.

We won't use this definition directly, but it's good to know.

Easy to prove equivalency to previous definition (see Lem. 3.4 in Bubeck's book).

Theorem (GD convergence for β -smooth functions.)

Let f be a β smooth convex function and assume we have $\|\mathbf{x}^* - \mathbf{x}^{(0)}\|_2 \leq R$. If we run GD for T steps, we have:

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T}$$

Corollary: If $T = O\left(\frac{\beta R^2}{\epsilon}\right)$ we have $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon$.

Compare this to $T = \underline{\underline{O\left(\frac{G^2 R^2}{\epsilon^2}\right)}}$ without a smoothness assumption.

GUARANTEED PROGRESS

Previously learning rate/step size η depended on G . Now choose it based on β .

$$\underline{\mathbf{x}^{(t+1)}} \leftarrow \underline{\mathbf{x}^{(t)}} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$$

Progress per step of gradient descent:

$$-\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$$

$$1. \underbrace{[f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] - \nabla f(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})}_{-\left(-\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\right)^T \left(\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\right)} \leq \frac{\beta}{2} \underbrace{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}\|_2^2}_{\left\|\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\right\|_2^2}$$

$$2. [f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] + \frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{\beta}{2} \underbrace{\left\|\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\right\|_2^2}_{\frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2}$$

$$3. \underline{f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)})} \geq \underline{\frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2}$$

CONVERGENCE GUARANTEE

Once we have the bound from the previous page, proving a convergence result isn't hard, but not obvious. A concise proof can be found in Page 15 in [Garrigos and Gower's notes.](#))

Theorem (GD convergence for β -smooth functions.)

Let f be a β smooth convex function and assume we have $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$. If we run GD for T steps with $\eta = \frac{1}{\beta}$ we have:

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T}$$

Corollary: If $T = O\left(\frac{\beta R^2}{\epsilon}\right)$ we have $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon$.

Where did we use convexity in this proof?

Progress per step of gradient descent:

$$1. [f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] - \nabla f(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) \leq \frac{\beta}{2} \|\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}\|_2^2.$$

$$2. [f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})] + \frac{1}{\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{\beta}{2} \|\frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})\|_2^2.$$

$$3. \underline{f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)})} \geq \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2.$$

STATIONARY POINTS

Definition (Stationary point)

For a differentiable function f , a stationary point is any x with:

$$\underline{\nabla f(x)} = 0$$

local/global minima - local/global maxima - saddle points



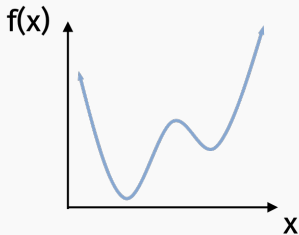
CONVERGENCE TO STATIONARY POINT

Theorem (Convergence to Stationary Point)

For any β -smooth differentiable function f (convex or not), if we run GD for T steps, we can find a point $\hat{\mathbf{x}}$ such that:

$$\underbrace{\|\nabla f(\hat{\mathbf{x}})\|_2^2}_{\leq} \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

Corollary: If $T \geq \frac{2\beta}{\epsilon}$, then $\underbrace{\|\nabla f(\hat{\mathbf{x}})\|_2^2}_{\leq} \leq \underbrace{\epsilon (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))}_{\leq}$.



TELESCOPING SUM PROOF

Theorem (Convergence to Stationary Point)

For any β -smooth differentiable function f (convex or not), if we run GD for T steps, we can find a point $\hat{\mathbf{x}}$ such that:

$$\|\nabla f(\hat{\mathbf{x}})\|_2^2 \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

$$\begin{aligned} &+ (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(1)})) + (f(\mathbf{x}^{(1)}) - f(\mathbf{x}^{(2)})) \\ &\vdots \\ &+ (f(\mathbf{x}^{(T-1)}) - f(\mathbf{x}^{(T)})) \end{aligned}$$

We have that $\frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \underline{f(\mathbf{x}^{(t)})} - \underline{f(\mathbf{x}^{(t+1)})}$. So:

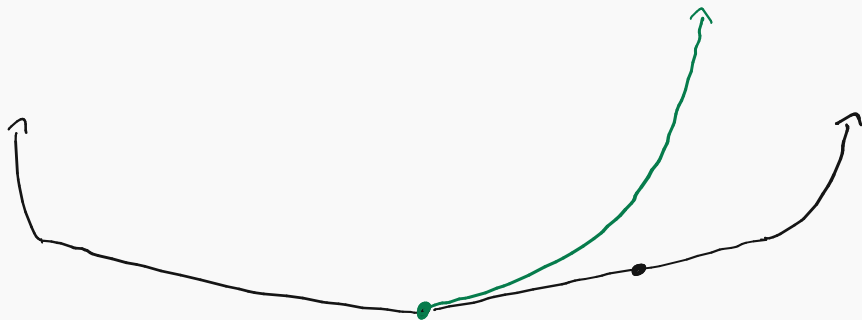
$$\sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \underline{f(\mathbf{x}^{(0)})} - \underline{f(\mathbf{x}^{(T)})}$$

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

$$\min_t \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{2\beta}{T} (f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*))$$

BACK TO CONVEX FUNCTIONS

I said it was a bit tricky to prove that $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2\beta R^2}{T}$ for convex functions. But we just easily proved that $\|\nabla f(\hat{\mathbf{x}})\|_2^2$ is small. Why doesn't this show we are close to the minimum?



Definition (α -strongly convex)

A convex function f is α -strongly convex if, for all \mathbf{x}, \mathbf{y}

$$\underbrace{[f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})}_{\geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2}$$

Compare to smoothness condition.

$$[f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2.$$

For a twice-differentiable scalar function f , equivalent to $f''(x) \geq \alpha$.

When f is convex, we always have that $f''(x) \geq 0$, so larger values of α correspond to a “stronger” condition.

Gradient descent for strongly convex functions:

- Choose number of steps T .
- For $i = 0, \dots, T$:
 - $\eta = \frac{2}{\alpha \cdot (i+1)}$
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)})$.

Theorem (GD convergence for α -strongly convex functions.)

Let f be an α -strongly convex function and assume we have that, for all \mathbf{x} , $\|\nabla f(\mathbf{x})\|_2 \leq G$. If we run GD for T steps (with adaptive step sizes) we have:

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2G^2}{\alpha T}$$

Corollary: If $T = O\left(\frac{G^2}{\alpha\epsilon}\right)$ we have $f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \epsilon$

CONVERGENCE GUARANTEE

$$f(\hat{x}) - f(x^*) \leq \epsilon$$

We could also have that f is both β -smooth and α -strongly convex.

Theorem (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\frac{2}{\beta} (f(x^{(T)}) - f(x^*)) \leq \underbrace{\|x^{(T)} - x^*\|_2^2}_{\leq e^{-T \frac{\alpha}{\beta}} \|x^{(0)} - x^*\|_2^2} \leq e^{-T \frac{\alpha}{\beta}} R^2$$

$\kappa = \frac{\beta}{\alpha}$ is called the “condition number” of f .

Is it better if κ is large or small?

$$e^{-T \frac{\alpha}{\beta}} = \epsilon$$

$$T = \frac{\beta}{\alpha} \log(1/\epsilon)$$

Converting to more familiar form: Using that fact the

$\nabla f(\mathbf{x}^*) = \mathbf{0}$ along with $\mathbf{x} = \mathbf{x}^*$ $\mathbf{y} = \mathbf{x}^{(T)}$

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq [f(\mathbf{y}) - f(\mathbf{x})] - \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

we have:

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \geq \frac{2}{\beta} [f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*)].$$

We also assume

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 \leq R^2.$$

CONVERGENCE GUARANTEE

Corollary (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\underline{f(\mathbf{x}^{(T)})} - \underline{f(\mathbf{x}^*)} \leq \frac{\cancel{\beta}}{2} e^{-T \frac{\alpha}{\beta}} \cdot \underline{R^2}$$

Corollary: If $T = O\left(\frac{\beta}{\alpha} \log(R\beta/\epsilon)\right)$ we have:

$$\underline{f(\mathbf{x}^{(T)})} - f(\mathbf{x}^*) \leq \epsilon$$

Only depend on $\log(1/\epsilon)$ instead of on $1/\epsilon$ or $1/\epsilon^2$!

χ^2

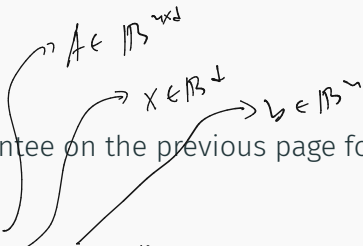
$2 \times$

2

SMOOTH, STRONGLY CONVEX OPTIMIZATION

We are going to prove the guarantee on the previous page for the special case of:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$



Goal: Get some of the key ideas across, introduces important concepts like the Hessian, and show the connection between conditioning and linear algebra.

THE HESSIAN

Let f be a twice differentiable function from $\underline{\mathbb{R}^d} \rightarrow \underline{\mathbb{R}}$. Let the **Hessian** $\underline{\nabla^2 f(x)}$ contain all of its second derivatives at a point x . So the Hessian is a $d \times d$ matrix and we have:



For vector x, v :

$$[\nabla^2 f(x)]_{j,k} = \frac{\partial^2 f}{\partial x_j \partial x_k}.$$

$$f(x + tv) \approx f(x) + t \langle v, \nabla f(x) \rangle$$

$$\nabla f(x + tv) \approx \nabla f(x) + t \underbrace{[\nabla^2 f(x)] v}_{d \times d} \quad \text{with } \nabla^2 f(x) \text{ being } d \times d$$

$$\lim_{t \rightarrow 0} \frac{\nabla f(x + tv) - \nabla f(x)}{t} = [\nabla^2 f(x)] v.$$

THE HESSIAN

Let f be a twice differentiable function from $\mathbb{R}^d \rightarrow \mathbb{R}$. Let the **Hessian** $\nabla^2 f(\mathbf{x})$ contain all of its second derivatives at a point \mathbf{x} . So the Hessian is a $d \times d$ matrix and we have:

$$\overset{x^2}{\nearrow} \quad [\nabla^2 f(\mathbf{x})]_{j,k} = \frac{\partial^2 f}{\partial x_j \partial x_k} \quad \swarrow \text{ith row of } \mathbf{a}.$$

Example: $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}^T \mathbf{a}^{(i)} - b^{(i)})^2$

$$\frac{\partial f}{\partial x_k} = \frac{1}{2} \sum_{i=1}^n 2 (\mathbf{x}^T \mathbf{a}^{(i)} - b^{(i)}) \cdot a_k^{(i)}$$

$$\frac{\partial^2 f}{\partial x_j \partial x_k} = \left(\sum_{i=1}^n a_j^{(i)} a_k^{(i)} \right) = a_j^T a_k$$

$$\nabla^2 f(\mathbf{x}) = \underline{\underline{A^T A}}$$

where a_j, a_k
are j th and k th
columns of
 A .

ALTERNATIVE DERIVATION

$$\underline{f(x)} = \frac{1}{2} \|Ax - b\|_2^2. \text{ Recall that } \nabla f(x) = \frac{1}{2} \cdot 2A^T(Ax - b) = A^T(Ax - b)$$

$$\lim_{t \rightarrow 0} \frac{\nabla f(x + tv) - \nabla f(x)}{t} = \underbrace{[\nabla^2 f(x)]}_{\text{Hessian}} v.$$

$$\frac{A^T(A(x+tv) - b) - A^T(Ax - b)}{t} = \frac{\cancel{A^T Ax} + A^T A tv - \cancel{A^T b} - \cancel{A^T Ax} + \cancel{A^T b}}{t}$$

$$\frac{\cancel{A^T A} tv}{\cancel{A}} = A^T A v$$

$$\downarrow$$

$$= \nabla^2 f(x)$$

IMPORTANT NOTE

The Hessian matrix is symmetric if for all j, k ,

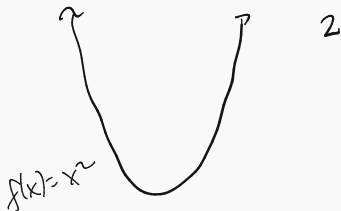
$$\frac{\partial^2 f}{\partial x_j \partial x_k} = \frac{\partial^2 f}{\partial x_k \partial x_j}$$

I.e. the order of differentiation does not matter. This is true whenever the second derivatives are continuous, which we will assume is the case.

CONVEXITY IN ONE DIMENSION

A twice-differentiable function $f: \mathbb{R} \rightarrow \mathbb{R}$ is :

- convex if and only if $f''(x) \geq 0$ for all x .
- β -smooth if $f''(x) \leq \beta$.
- α -strongly convex if $f''(x) \geq \alpha$.



How do these statements generalize to the case when f has a vector input, so the second derivative is a matrix $\nabla^2 f(\mathbf{x})$?

Claim: If f is twice differentiable, then it is convex if and only if the matrix $\nabla^2 f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} .

Definition (Positive Semidefinite (PSD))

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is positive semidefinite (PSD) for any vector $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{y}^T \mathbf{H} \mathbf{y} \geq 0$.

This is a natural notion of “positivity” for symmetric matrices. To denote that \mathbf{H} is PSD we will typically use “Loewner order” notation (`\succeq` in LaTeX):

$$\mathbf{H} \succeq 0.$$

We write $\mathbf{B} \succeq \mathbf{A}$ or equivalently $\mathbf{A} \preceq \mathbf{B}$ to denote that $(\mathbf{B} - \mathbf{A})$ is positive semidefinite. This gives a partial ordering on matrices.

$$\mathbf{B} - \mathbf{A} \succeq 0 \quad \mathbf{A} \preceq \mathbf{B} \quad \mathbf{B} - \mathbf{A} \succeq 0.$$

HESSIAN MATRICES AND POSITIVE SEMIDEFINITENESS

Claim: If f is twice differentiable, then it is convex if and only if the matrix $\nabla^2 f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} .

Definition (Positive Semidefinite (PSD))

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is positive semidefinite (PSD) for any vector $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{y}^T \mathbf{H} \mathbf{y} \geq 0$.

For the least squares regression loss function:

$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, $\nabla^2 f(\mathbf{x}) = \mathbf{A}^T \mathbf{A}$ for all \mathbf{x} . Is $\nabla^2 f(\mathbf{x})$ PSD?

For any \mathbf{y} , $\mathbf{y}^T \mathbf{A}^T \mathbf{A} \mathbf{y} \geq 0$.

$\langle \mathbf{A}\mathbf{y}, \mathbf{A}\mathbf{y} \rangle = \|\mathbf{A}\mathbf{y}\|_2^2 \geq 0$.

If f is β -smooth and α -strongly convex then at any point \mathbf{x} , the Hessian $\nabla^2 f(\mathbf{x})$ satisfies:

$$\underline{\alpha \mathbf{I}} \preceq \underline{\nabla^2 f(\mathbf{x})} \preceq \underline{\beta \mathbf{I}},$$

where \mathbf{I} is a $d \times d$ identity matrix.

This is the natural matrix generalization of the statement for scalar valued functions:

$$\alpha \leq f'(x) \leq \beta.$$

SMOOTH AND STRONGLY CONVEX HESSIAN

$$\alpha I_{d \times d} \preceq \nabla^2 f(x) \preceq \beta I_{d \times d}.$$

Equivalently for any z ,

$$\alpha \|z\|_2^2 \leq z^T [\nabla^2 f(x)] z \leq \beta \|z\|_2^2.$$

$$\beta I - \nabla^2 f(x) \succeq 0$$

Is PSD.

For all z ,

$$z^T (\beta I - \nabla^2 f(x)) z \geq 0$$
$$\beta z^T z - z^T \nabla^2 f(x) z \geq 0$$

$\underbrace{\beta z^T z}_{\|z\|_2^2}$

SIMPLE EXAMPLE

Let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ where \mathbf{D} is a diagonal matrix. For now imagine we're in two dimensions: $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$, $\mathbf{D} = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$.

What are α, β for this problem?

$$\nearrow \mathbf{D}^T \mathbf{D} = \mathbf{D}^2 = \begin{pmatrix} d_1^2 & 0 \\ 0 & d_2^2 \end{pmatrix}$$

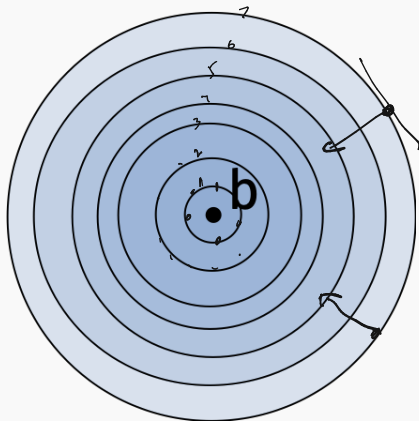
$$\alpha \|\mathbf{z}\|_2^2 \leq \mathbf{z}^T [\nabla^2 f(\mathbf{x})] \mathbf{z} \leq \beta \|\mathbf{z}\|_2^2$$

$$\swarrow \min(d_1^2, d_2^2)$$

$$\downarrow \max(d_1^2, d_2^2)$$

$$\mathbf{z} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \mathbf{z} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

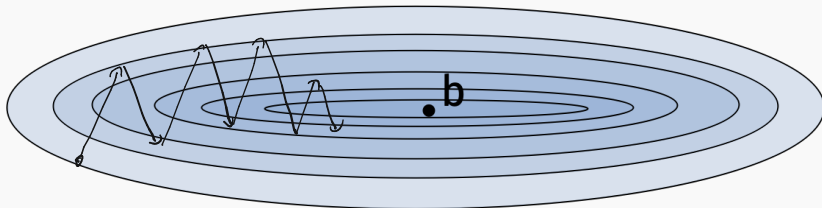
GEOMETRIC VIEW



Level sets of $\frac{1}{2} \| \mathbf{Dx} - \mathbf{b} \|_2^2$ when $\underline{d_1^2} = 1, \underline{d_2^2} = 1$.

$$= \frac{1}{2} \| \mathbf{x} - \mathbf{b} \|_2^2$$

GEOMETRIC VIEW

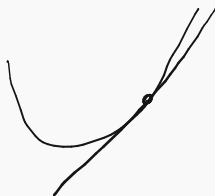


Level sets of $\frac{1}{2} \|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ when $d_1^2 = \frac{1}{3}$, $d_2^2 = 2$.

$$\kappa = 2$$
$$\alpha = 1/3$$

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{3}{0.5} = 6$$

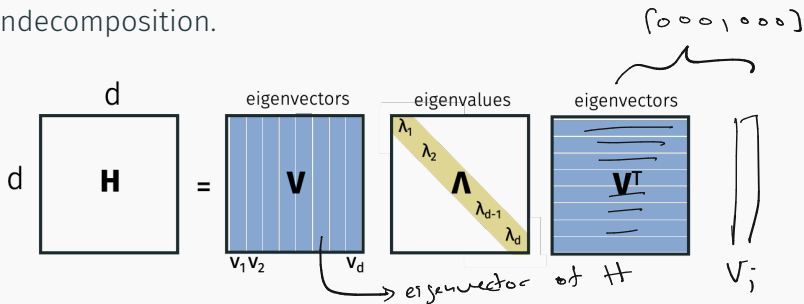
condition #



What about non-diagonal D?

EIGENDECOMPOSITION VIEW

Any symmetric matrix \mathbf{H} has an orthogonal, real valued eigendecomposition.



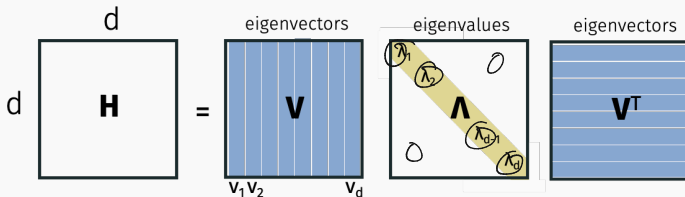
Here \mathbf{V} is square and orthogonal, so $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}$. And for each \mathbf{v}_i , we have:

$$\underline{\mathbf{H}} \underline{\mathbf{v}}_j = \lambda_j \underline{\mathbf{v}}_j.$$

By definition, that's what makes $\mathbf{v}_1, \dots, \mathbf{v}_d$ eigenvectors.

EIGENDECOMPOSITION VIEW

Recall $\mathbf{V}\mathbf{V}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$.



Claim: \mathbf{H} is PSD $\Leftrightarrow \lambda_1, \dots, \lambda_d \geq 0$. **Proof for \Leftarrow :**

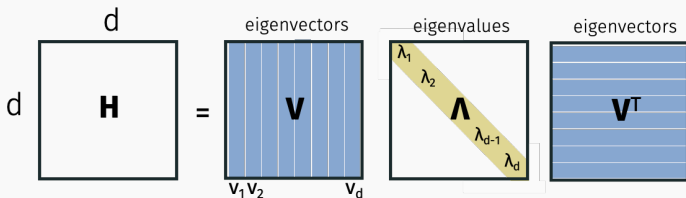
$$\mathbf{H} = \mathbf{V} \sqrt{\mathbf{\Lambda}} \sqrt{\mathbf{\Lambda}} \mathbf{V}^T = (\sqrt{\mathbf{\Lambda}} \mathbf{V}^T)^T (\sqrt{\mathbf{\Lambda}} \mathbf{V}^T) = \underline{\underline{\mathbf{C}^T \mathbf{C}}}$$

$$\mathbf{z}^T \mathbf{H} \mathbf{z} = \|\mathbf{C} \mathbf{z}\|_2^2 \geq 0.$$

EIGENDECOMPOSITION VIEW

Recall $\mathbf{W}\mathbf{W}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$.

$$\lambda_1 \geq \lambda_2 \dots \geq \lambda_d$$



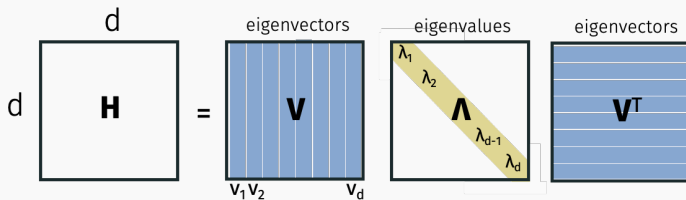
Claim: $\alpha \mathbf{I} \preceq \mathbf{H} \preceq \beta \mathbf{I} \Leftrightarrow \alpha \leq \lambda_d \leq \dots \leq \lambda_1 \leq \beta$.

$$\mathbf{H} \preceq \beta \mathbf{I} \quad (\beta \mathbf{I} - \mathbf{H}) \succeq \mathbf{0}$$

$$\beta \cdot \mathbf{V}\mathbf{V}^T - \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}(\beta \mathbf{I} - \mathbf{\Lambda})\mathbf{V}^T$$

$$(\beta \mathbf{I} - \mathbf{\Lambda})_{ii} \geq 0 \quad \text{for all } i.$$

EIGENDECOMPOSITION VIEW



Recall that if $\lambda_{\max}(\mathbf{H})$ and $\lambda_{\min}(\mathbf{H})$ be the smallest and largest eigenvalues of \mathbf{H} , then for all \mathbf{z} we have:

$$\mathbf{z}^T \mathbf{H} \mathbf{z} \leq \lambda_{\max}(\mathbf{H}) \cdot \|\mathbf{z}\|^2$$

$$\mathbf{z}^T \mathbf{H} \mathbf{z} \geq \lambda_{\min}(\mathbf{H}) \cdot \|\mathbf{z}\|^2$$

EIGENDECOMPOSITION VIEW

If for all \mathbf{x} the maximum eigenvalue of $\nabla^2 f(\mathbf{x})$ is $\leq \underline{\beta}$ and the minimum eigenvalue of $\nabla^2 f(\mathbf{x})$ is $\geq \alpha$ then $f(\mathbf{x})$ is β -smooth and α -strongly convex.

Note that for $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$, we have that, for all \mathbf{x} , $\nabla^2 f(\mathbf{x}) = \underline{\mathbf{A}^T \mathbf{A}}$. So, we can take $\alpha = \underline{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}$ and $\beta = \underline{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$.

Theorem (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for S steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\|x^{(S)} - x^*\|_2 \leq e^{-S/\kappa} \|x^{(0)} - x^*\|_2$$

Goal: Prove for $f(x) = \frac{1}{2} \|Ax - b\|_2^2$.

Let $\lambda_{\max} = \lambda_{\max}(A^T A)$ and set step size $\eta = \frac{1}{\lambda_{\max}}$. Gradient descent update is:

$$\left(\underline{x^{(t+1)}} = \underline{x^{(t)}} - \frac{1}{\lambda_{\max}} A^T (Ax^{(t)} - b) \right)$$

\swarrow
 μ

ALTERNATIVE VIEW OF GRADIENT DESCENT

Richardson Iteration view:

$$\nabla f(x^*) = 0$$

$$A^T(Ax^* - b) = 0$$

$$\underline{A^T A x^* = A^T b}$$

$$\underline{(x^{(t+1)} - x^*)} = \left(I - \frac{1}{\lambda_{\max}} A^T A \right) \underline{(x^{(t)} - x^*)}$$

$$x^{(t)} - x^* - \frac{1}{\lambda_{\max}} A^T A x^{(t)} + \frac{1}{\lambda_{\max}} A^T A x^*$$

$$= x^{(t)} - x^* - \frac{1}{\lambda_{\max}} A^T A x^{(t)} + \frac{1}{\lambda_{\max}} A^T b$$

$$\frac{1}{\lambda_{\max}} A^T (A x^{(t)} - b)$$

$$= x^{(t)} - \frac{1}{\lambda_{\max}} A^T (A x^{(t)} - b) - x^*$$

$$x^{(t+1)}$$

UNROLLED GRADIENT DESCENT

$$(x^{(s)} - x^*) = \left(I - \frac{1}{\lambda_{\max}} A^T A \right)^s (x^{(0)} - x^*)$$

$$x^{(1)} - x^* = \left(I - \frac{1}{\lambda_{\max}} A^T A \right) (x^{(0)} - x^*)$$

$$(x^{(2)} - x^*) = \left(I - \frac{1}{\lambda_{\max}} A^T A \right) (x^{(1)} - x^*)$$

⋮

$$(x^{(s)} - x^*) = \left(I - \frac{1}{\lambda_{\max}} A^T A \right)^s (x^{(0)} - x^*)$$

UNROLLED GRADIENT DESCENT

$$(x^{(S)} - x^*) = \left(\underbrace{\left(I - \frac{1}{\lambda_{\max}} A^T A \right)^S}_{\mathcal{B}} (x^{(0)} - x^*) \right)$$

Conclusion: $\|x^{(S)} - x^*\|_2 \leq (x^{(0)} - x^*)^T \mathcal{B}^S (x^{(0)} - x^*)$

$$\leq \|x^{(0)} - x^*\|_2 \cdot \lambda_{\max}(\mathcal{B}^S).$$

H

$$\|z^T H z\| \leq \|z\|_2^2 (\lambda_{\max}(H))$$

Approach: Show that the maximum eigenvalue of $\left(I - \frac{1}{\lambda_{\max}} A^T A \right)^{2S}$ is small – i.e., bounded by $\underline{\underline{e^{-S/\kappa}}} = \epsilon$.

UNROLLED GRADIENT DESCENT

$$(x^{(S)} - x^*) = \left(I - \frac{1}{\lambda_{\max}} A^T A \right)^S (x^{(0)} - x^*)$$

What is the maximum eigenvalue of the symmetric matrix $\left(I - \frac{1}{\lambda_{\max}} A^T A \right)$ in terms of the eigenvalues of $A^T A$?

$$V V^T - \frac{1}{\lambda_{\max}} V \Lambda V^T = V \left(I - \frac{1}{\lambda_{\max}} \Lambda \right) V^T = I - \frac{1}{\lambda_{\max}} (A^T A)$$

$$I - \frac{\lambda_1}{\lambda_{\max}} \quad \lambda_{\max} \left(I - \frac{1}{\lambda_{\max}} A^T A \right) = I - \frac{\lambda_{\max}(A^T A)}{\lambda_{\max}(A^T A)}$$

$$I - \frac{\lambda_2}{\lambda_{\max}}$$

$$\vdots$$

$$I - \frac{\lambda_d = \lambda_{\max}}{\lambda_{\max}}$$

UNROLLED GRADIENT DESCENT

$$\left(1 - \frac{1}{\kappa}\right)^s \leq \frac{1}{e} \quad \text{for } \kappa \geq 2$$

$$(\mathbf{x}^{(s)} - \mathbf{x}^*) = \left(\mathbf{I} - \frac{1}{\lambda_{\max}} \mathbf{A}^T \mathbf{A} \right)^s (\mathbf{x}^{(0)} - \mathbf{x}^*)$$

What is the maximum eigenvalue of $\underbrace{\left(\mathbf{I} - \frac{1}{\lambda_{\max}} \mathbf{A}^T \mathbf{A} \right)}_{\mathbf{B}}$ ^(2S)?

$$\mathbf{B} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

$$\mathbf{B}^s = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \dots \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T = \mathbf{V} \mathbf{\Lambda}^s \mathbf{V}^T$$

$$\lambda_{\max}(\mathbf{B}) = 1 - \frac{\lambda_{\min}}{\lambda_{\max}} \quad \lambda_{\max}(\mathbf{B}^{2s}) = \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}} \right)^{2s}$$

$$= \left(\left(1 - \frac{1}{\kappa} \right)^{2s\kappa} \right)^{1/\kappa} \leq \left(\frac{1}{e} \right)^{s/\kappa} = e^{-s/\kappa}$$

ACCELERATION

ACCELERATED GRADIENT DESCENT

Nesterov's accelerated gradient descent:

- $\mathbf{x}^{(0)} = \mathbf{y}^{(1)} = \mathbf{z}^{(1)}$

- For $t = 1, \dots, T$

- $\mathbf{y}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$

- $\mathbf{x}^{(t+1)} = \left(1 + \frac{\sqrt{\kappa}-1}{\sqrt{\kappa+1}}\right) \mathbf{y}^{(t+1)} + \frac{\sqrt{\kappa}-1}{\sqrt{\kappa+1}} (\mathbf{y}^{(t+1)} - \mathbf{y}^{(t)})$

$$e^{-S/\kappa}$$
$$S = \kappa \log(1/\epsilon)$$

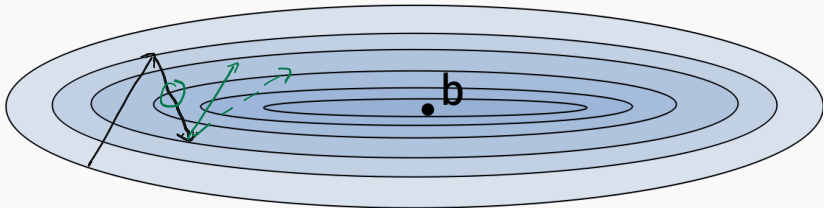
Theorem (AGD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run AGD for S steps we have:

$$\|\mathbf{x}^{(S)} - \mathbf{x}^*\|_2 \leq e^{-S/\sqrt{\kappa}} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$$

Corollary: If $T = O(\sqrt{\kappa} \log(\beta R/\epsilon))$ achieve error ϵ .

INTUITION BEHIND ACCELERATION



Level sets of $\|Ax - b\|_2^2$.

Other terms for similar ideas:

- Momentum
- Heavy-ball methods

What if we look back beyond two iterates?

BREAK

Second part of class:

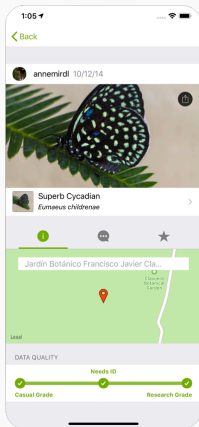
- Basics of Online Learning + Optimization.
- Introduction to Regret Analysis.
- Application to analyzing Stochastic Gradient Descent.

Many machine learning problems are solved in an online setting with constantly changing data.

- Spam filters are incrementally updated and adapt as they see more examples of spam over time.
- Image classification systems learn from mistakes over time (often based on user feedback).
- Content recommendation systems adapt to user behavior and clicks (which may not be a good thing...)

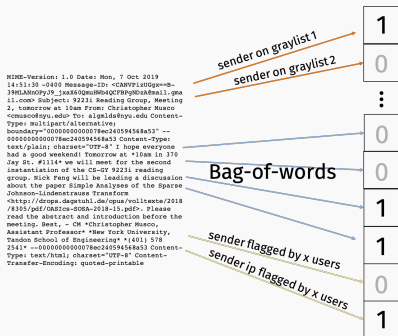
Plant identification via iNaturalist app.

(California Academy of Science + National Geographic)



- When the app fails, image is classified via crowdsourcing (backed by huge network of amateurs and experts).
- Single model that is updated constantly, not retrained in batches.

Machine learning based email spam filtering.



Markers for spam change overtime, so model might change.

Machine learning based email spam filtering.



Markers for spam change overtime, so model might change.

Choose some model $M_{\mathbf{x}}$ parameterized by parameters \mathbf{x} and some loss function ℓ . At time steps $1, \dots, T$, receive data vectors $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(T)}$.

- At each time step, we pick (“play”) a parameter vector $\mathbf{x}^{(i)}$.
- Make prediction $\tilde{y}^{(i)} = M_{\mathbf{x}^{(i)}}(\mathbf{a}_i)$.
- Then told true value or label $y^{(i)}$. Possibly use this information to choose a new $\mathbf{x}^{(i+1)}$.
- Goal is to minimize cumulative loss:

$$L = \sum_{i=1}^n \ell(\mathbf{x}^{(i)}, \mathbf{a}^{(i)}, y^{(i)})$$

For example, for a regression problem we might use the ℓ_2 loss:

$$\ell(\mathbf{x}^{(i)}, \mathbf{a}^{(i)}, y^{(i)}) = \left| \langle \mathbf{x}^{(i)}, \mathbf{a}^{(i)} \rangle - y^{(i)} \right|^2.$$

For classification, we could use logistic/cross-entropy loss.

Abstraction as optimization problem: Instead of a single objective function f , we have a single (initially unknown) function $f_1, \dots, f_T : \mathbb{R}^d \rightarrow \mathbb{R}$ for each time step.

- For time step $i \in 1, \dots, T$, select vector $\mathbf{x}^{(i)}$.
- Observe f_i and pay cost $f_i(\mathbf{x}^{(i)})$
- Goal is to minimize $\sum_{i=1}^T f_i(\mathbf{x}^{(i)})$.

We make no assumptions that f_1, \dots, f_T are related to each other at all!

In offline optimization, we wanted to find $\hat{\mathbf{x}}$ satisfying $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x})$. Ask for a similar thing here.

Objective: Choose $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ so that:

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right] + \epsilon.$$

Here ϵ is called the **regret** of our solution sequence $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(T)}$.

We typically ϵ to be growing sublinearly in T .

Regret compares to the best fixed solution in hindsight.

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right] + \epsilon.$$

It's very possible that $\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) < \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right]$. Could we hope for something stronger?

Exercise: Argue that the following is impossible to achieve:

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\sum_{i=1}^T \min_{\mathbf{x}} f_i(\mathbf{x}) \right] + \epsilon.$$

Convex functions:

$$f_1(x) = |x - h_1|$$

$$\vdots$$

$$f_n(x) = |x - h_T|$$

where h_1, \dots, h_T are i.i.d. uniform $\{0, 1\}$.

$$\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \leq \left[\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x}) \right] + \epsilon.$$

Beautiful balance:

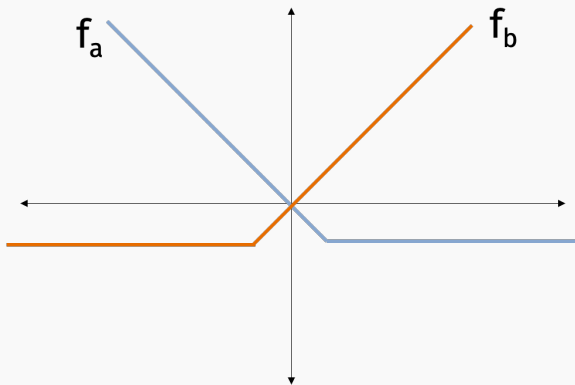
- Either f_1, \dots, f_T are similar or changing slowly, so we can learn predict f_i from earlier functions.
- Or f_1, \dots, f_T are very different, in which case $\min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x})$ is large, so regret bound is easy to achieve.
- Or we live somewhere in the middle.

Follow-the-leader algorithm:

- Choose $\mathbf{x}^{(0)}$.
- For $i = 1, \dots, T$:
 - Let $\mathbf{x}^{(i)} = \arg \min_{\mathbf{x}} \sum_{j=1}^{i-1} f_j(\mathbf{x})$.
 - Play $\mathbf{x}^{(i)}$.
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.

Simple and intuitive, but there are two issues with this approach. One is computational, one is related to the accuracy.

Hard case:



Online Gradient descent:

- Choose $\mathbf{x}^{(1)}$ and $\eta = \frac{R}{G\sqrt{T}}$.
- For $i = 1, \dots, T$:
 - Play $\mathbf{x}^{(i)}$.
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_i(\mathbf{x}^{(i)})$

If $f_1, \dots, f_T = f$ are all the same, this looks a lot like regular gradient descent. We update parameters using the gradient ∇f at each step.

ONLINE GRADIENT DESCENT (OGD)

$\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x})$ (the offline optimum)

Assume:

- f_1, \dots, f_T are all convex.
- Each is G -Lipschitz: for all \mathbf{x}, i , $\|\nabla f_i(\mathbf{x})\|_2 \leq G$.
- Starting radius: $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$.

Online Gradient descent:

- Choose $\mathbf{x}^{(1)}$ and $\eta = \frac{R}{G\sqrt{T}}$.
- For $i = 1, \dots, T$:
 - Play $\mathbf{x}^{(i)}$.
 - Observe f_i and incur cost $f_i(\mathbf{x}^{(i)})$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_i(\mathbf{x}^{(i)})$

Let $\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{i=1}^T f_i(\mathbf{x})$ (the offline optimum)

Theorem (OGD Regret Bound)

After T steps, $\epsilon = \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right] \leq RG\sqrt{T}$.

Average regret overtime is bounded by $\frac{\epsilon}{T} \leq \frac{RG}{\sqrt{T}}$.

Goes $\rightarrow 0$ as $T \rightarrow \infty$.

All this with no assumptions on how f_1, \dots, f_T relate to each other! They could have even been chosen **adversarially** – e.g. with f_i depending on our choice of \mathbf{x}_i and all previous choices.

Theorem (OGD Regret Bound)

After T steps, $\epsilon = \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right] \leq RG\sqrt{T}$.

Claim 1: For all $i = 1, \dots, T$,

$$f_i(\mathbf{x}^{(i)}) - f_i(\mathbf{x}^*) \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

(Same proof for standard GD. Only uses convexity of f_i .)

Theorem (OGD Regret Bound)

After T steps, $\epsilon = \left[\sum_{i=1}^T f_i(\mathbf{x}^{(i)}) \right] - \left[\sum_{i=1}^T f_i(\mathbf{x}^*) \right] \leq RG\sqrt{T}$.

Claim 1: For all $i = 1, \dots, T$,

$$f_i(\mathbf{x}^{(i)}) - f_i(\mathbf{x}^*) \leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Telescoping Sum:

$$\begin{aligned} \sum_{i=1}^T \left[f_i(\mathbf{x}^{(i)}) - f_i(\mathbf{x}^*) \right] &\leq \frac{\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2} \\ &\leq \frac{R^2}{2\eta} + \frac{T\eta G^2}{2} \end{aligned}$$

STOCHASTIC GRADIENT DESCENT (SGD)

Efficient offline optimization method for functions f with finite sum structure:

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x}).$$

Goal is to find $\hat{\mathbf{x}}$ such that $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$.

- The most widely use optimization algorithm in modern machine learning.
- Easily analyzed as a special case of online gradient descent!

Recall the machine learning setup. In empirical risk minimization, we can typically write:

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$$

where f_i is the loss function for a particular data example $(\mathbf{a}^{(i)}, y^{(i)})$.

Example: least squares linear regression.

$$f(\mathbf{x}) = \sum_{i=1}^n (\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)})^2$$

Note that by linearity, $\nabla f(\mathbf{x}) = \sum_{i=1}^n \nabla f_i(\mathbf{x})$.

Main idea: Use random approximate gradient in place of actual gradient.

Pick random $j \in 1, \dots, n$ and update \mathbf{x} using $\nabla f_j(\mathbf{x})$.

$$\mathbb{E} [\nabla f_j(\mathbf{x})] = \frac{1}{n} \nabla f(\mathbf{x}).$$

$n\nabla f_j(\mathbf{x})$ is an unbiased estimate for the true gradient $\nabla f(\mathbf{x})$, but can often be computed in a $1/n$ fraction of the time!

Trade slower convergence for cheaper iterations.

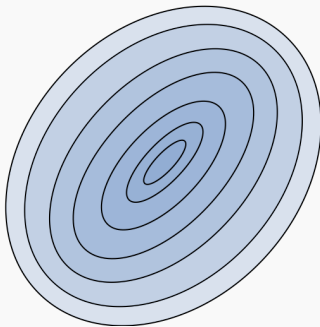
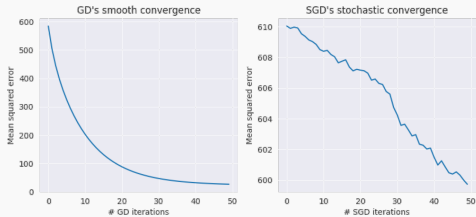
Stochastic first-order oracle for $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$.

- **Function Query:** For any chosen j, \mathbf{x} , return $f_j(\mathbf{x})$
- **Gradient Query:** For any chosen j, \mathbf{x} , return $\nabla f_j(\mathbf{x})$

Stochastic Gradient descent:

- Choose starting vector $\mathbf{x}^{(1)}$, step size η
- For $i = 1, \dots, T$:
 - Pick random $j_i \in 1, \dots, n$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^T \mathbf{x}^{(i)}$

VISUALIZING SGD



STOCHASTIC GRADIENT DESCENT

Assume:

- Finite sum structure: $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, with f_1, \dots, f_n all convex.
- Lipschitz functions: for all \mathbf{x}, j , $\|\nabla f_j(\mathbf{x})\|_2 \leq \frac{G'}{n}$.
 - What does this imply about Lipschitz constant of f ?
- Starting radius: $\|\mathbf{x}^* - \mathbf{x}^{(1)}\|_2 \leq R$.

Stochastic Gradient descent:

- Choose $\mathbf{x}^{(1)}$, steps T , step size $\eta = \frac{R}{G'\sqrt{T}}$.
- For $i = 1, \dots, T$:
 - Pick random $j_i \in 1, \dots, n$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f_{j_i}(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^T \mathbf{x}^{(i)}$

Approach: View as online gradient descent run on function sequence f_{j_1}, \dots, f_{j_T} .

Only use the fact that step equals gradient in expectation.

JENSEN'S INEQUALITY

For a convex function f and points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$

$$f\left(\frac{1}{t} \cdot \mathbf{x}^{(1)} + \dots + \frac{1}{t} \cdot \mathbf{x}^{(t)}\right) \leq \frac{1}{t} \cdot f(\mathbf{x}^{(1)}) + \dots + \frac{1}{t} \cdot f(\mathbf{x}^{(t)})$$

Claim (SGD Convergence)

After $T = \frac{R^2 G^2}{\epsilon^2}$ iterations:

$$\mathbb{E} [f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

Claim 1:

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{i=1}^T [f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)]$$

Prove using Jensen's Inequality:

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\begin{aligned}\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] &\leq \frac{1}{T} \sum_{i=1}^T \mathbb{E}[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \\ &= \frac{1}{T} \sum_{i=1}^T n \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)]\end{aligned}$$

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\begin{aligned}\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] &\leq \frac{1}{T} \sum_{i=1}^T \mathbb{E}[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \\ &= \frac{1}{T} \sum_{i=1}^T n \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)] \\ &\leq \frac{n}{T} \cdot \mathbb{E} \left[\sum_{i=1}^T f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^{offline}) \right],\end{aligned}$$

where $\mathbf{x}^{offline} = \arg \min_{\mathbf{x}} \sum_{i=1}^T f_{j_i}(\mathbf{x})$.

Claim (SGD Convergence)

After $T = \frac{R^2 G'^2}{\epsilon^2}$ iterations:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon.$$

$$\begin{aligned}\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] &\leq \frac{1}{T} \sum_{i=1}^T \mathbb{E}[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)] \\ &= \frac{1}{T} \sum_{i=1}^T n \mathbb{E}[f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^*)] \\ &\leq \frac{n}{T} \mathbb{E}\left[\sum_{i=1}^T f_{j_i}(\mathbf{x}^{(i)}) - f_{j_i}(\mathbf{x}^{\text{offline}})\right] \\ &\leq \frac{n}{T} \cdot \left(R \cdot \frac{G'}{n} \cdot \sqrt{T}\right) \quad (\text{by OGD guarantee.})\end{aligned}$$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

Number of iterations for error ϵ :

- Gradient Descent: $T = \frac{R^2 G^2}{\epsilon^2}$.
- Stochastic Gradient Descent: $T = \frac{R^2 G'^2}{\epsilon^2}$.

Always have $G \leq G'$:

$$\begin{aligned}\max_{\mathbf{x}} \|\nabla f(\mathbf{x})\|_2 &\leq \max_{\mathbf{x}} (\|\nabla f_1(\mathbf{x})\|_2 + \dots + \|\nabla f_n(\mathbf{x})\|_2) \\ &\leq \max_{\mathbf{x}} (\|\nabla f_1(\mathbf{x})\|_2) + \dots + \max_{\mathbf{x}} (\|\nabla f_n(\mathbf{x})\|_2) \\ &\leq n \cdot \frac{G'}{n} = G'.\end{aligned}$$

So GD converges strictly faster than *SGD*.

But for a fair comparison:

- SGD cost = (# of iterations) $\cdot O(1)$
- GD cost = (# of iterations) $\cdot O(n)$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

We always have $G \leq G'$. When it is much smaller then GD will perform better. When it is closer to this upper bound, SGD will perform better.

What is an extreme case where $G = G'$?

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

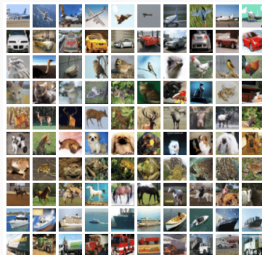
What if each gradient $\nabla f_i(\mathbf{x})$ looks like random vectors in \mathbb{R}^d ?
E.g. with $\mathcal{N}(0, 1)$ entries?

$$\mathbb{E} [\|\nabla f_i(\mathbf{x})\|_2^2] =$$

$$\mathbb{E} [\|\nabla f(\mathbf{x})\|_2^2] = \mathbb{E} \left[\left\| \sum_{i=1}^n \nabla f_i(\mathbf{x}) \right\|_2^2 \right] =$$

STOCHASTIC VS. FULL BATCH GRADIENT DESCENT

Takeaway: SGD performs better when there is more structure or repetition in the data set.



PRECONDITIONING

Main idea: Instead of minimizing $f(\mathbf{x})$, find another function $g(\mathbf{x})$ with the same minimum but which is better suited for first order optimization (e.g., has a smaller conditioner number).

Claim: Let $h(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an invertible function. Let $g(\mathbf{x}) = f(h(\mathbf{x}))$. Then

$$\min_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{x}} g(\mathbf{x}) \quad \text{and} \quad \arg \min_{\mathbf{x}} f(\mathbf{x}) = h \left(\arg \min_{\mathbf{x}} g(\mathbf{x}) \right).$$

First Goal: We need $g(\mathbf{x})$ to still be convex.

Claim: Let \mathbf{P} be an invertible $d \times d$ matrix and let $g(\mathbf{x}) = f(\mathbf{P}\mathbf{x})$.

$g(\mathbf{x})$ is always convex.

Second Goal:

$g(\mathbf{x})$ should have better condition number κ than $f(\mathbf{x})$.

Example:

- $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$. $\kappa_f = \frac{\lambda_1(\mathbf{A}^T\mathbf{A})}{\lambda_d(\mathbf{A}^T\mathbf{A})}$.
- $g(\mathbf{x}) = \|\mathbf{APx} - \mathbf{b}\|_2^2$. $\kappa_g = \frac{\lambda_1(\mathbf{P}^T\mathbf{A}^T\mathbf{AP})}{\lambda_d(\mathbf{P}^T\mathbf{A}^T\mathbf{AP})}$.

Third Goal: \mathbf{P} should be easy to compute.

Many, many problem specific preconditioners are used in practice. Their design is usually a heuristic process.

Example: Diagonal preconditioner.

- Let $\mathbf{D} = \text{diag}(\mathbf{A}^T\mathbf{A})$
- Intuitively, we roughly have that $\mathbf{D} \approx \mathbf{A}^T\mathbf{A}$.
- Let $\mathbf{P} = \sqrt{\mathbf{D}^{-1}}$

\mathbf{P} is often called a **Jacobi preconditioner**. Often works very well in practice!

DIAGONAL PRECONDITIONER

A =

| | | | | |
|-------|----|------|--------|-----|
| -734 | 1 | 33 | 9111 | 0 |
| -31 | -2 | 108 | 5946 | -19 |
| 232 | -1 | 101 | 3502 | 10 |
| 426 | 0 | -65 | 12503 | 9 |
| -373 | 0 | 26 | 9298 | 0 |
| -236 | -2 | -94 | 2398 | -1 |
| 2024 | 0 | -132 | -6904 | -25 |
| -2258 | -1 | 92 | -6516 | 6 |
| 2229 | 0 | 0 | 11921 | -22 |
| 338 | 1 | -5 | -16118 | -23 |

```
>> cond(A'*A)
```

```
ans =
```

```
8.4145e+07
```

```
>> P = sqrt(inv(diag(diag(A'*A))));
```

```
>> cond(P*A'*A*P)
```

```
ans =
```

```
10.3878
```

Another view: If $g(\mathbf{x}) = f(\mathbf{P}\mathbf{x})$ then $\nabla g(\mathbf{x}) = \mathbf{P}^T \nabla f(\mathbf{P}\mathbf{x})$.

$\nabla g(\mathbf{x}) = \mathbf{P} \nabla f(\mathbf{P}\mathbf{x})$ when \mathbf{P} is symmetric.

Gradient descent on g :

- For $t = 1, \dots, T$,
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \mathbf{P} [\nabla f(\mathbf{P}\mathbf{x}^{(t)})]$

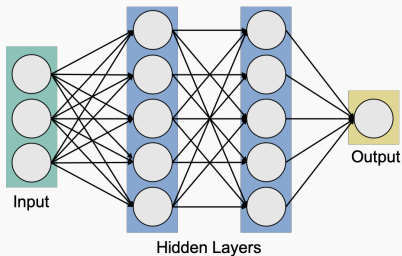
Gradient descent on g :

- For $t = 1, \dots, T$,
 - $\mathbf{y}^{(t+1)} = \mathbf{y}^{(t)} - \eta \mathbf{P}^2 [\nabla f(\mathbf{y}^{(t)})]$

When \mathbf{P} is diagonal, this is just gradient descent with a different step size for each parameter!

Algorithms based on this idea:

- AdaGrad
- RMSprop
- Adam optimizer



(Pretty much all of the most widely used optimization methods for training neural networks.)

COORDINATE DESCENT

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Gradient Descent: When $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, approximate $\nabla f(\mathbf{x})$ with $\nabla f_i(\mathbf{x})$ for randomly chosen i .

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Coordinate Descent: Only compute a single random entry of $\nabla f(\mathbf{x})$ on each iteration:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\mathbf{x}) \end{bmatrix} \qquad \nabla_{ij} f(\mathbf{x}) = \begin{bmatrix} 0 \\ \frac{\partial f}{\partial x_i}(\mathbf{x}) \\ \vdots \\ 0 \end{bmatrix}$$

Update: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta \nabla_{ij} f(\mathbf{x}^{(t)})$.