# CS-GY 6763: LECTURE 6
# GRADIENT DESCENT AND PROJECTED
# GRADIENT DESCENT

NYU Tandon School of Engineering, Prof. Rajesh Jayaram

## PROJECT

- HW Due this Friday 3/11 by end of day.
- If doing final project, start looking at papers, thinking about research problems (reach out to me if you need help).
- HW#3 released next week.
- Midterm during first half of class, 3/21
- Midterm prep sheet to be posted soon.

## NEW UNIT: CONTINUOUS OPTIMIZATION

Have some function $f : \mathbb{R}^d \to \mathbb{R}$. Want to find $\mathbf{x}^*$ such that:

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}).$$

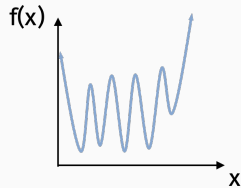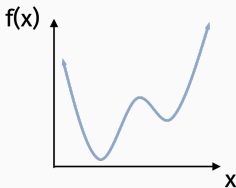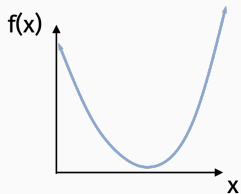Or at least $\hat{\mathbf{x}}$ which is close to a minimum. E.g.
$f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon$
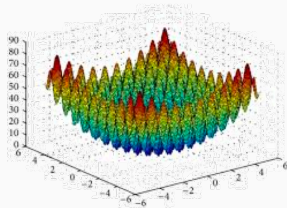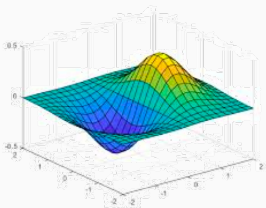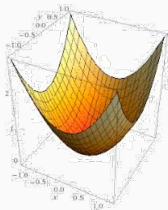
Often we have some additional constraints:

- $\mathbf{x} > 0$.
- $\|\mathbf{x}\|_2 \leq R$, $\|\mathbf{x}\|_1 \leq R$.
- $\mathbf{a}^T \mathbf{x} > c$.

**Dimension** $d = 1$:



**Dimension** $d = 2$:

## OPTIMIZATION IN MACHINE LEARNING

**Continuouos optimization is the foundation of modern machine learning.**

**Supervised learning:** Want to learn a model that maps <u>inputs</u>

- numerical data vectors
- images, video
- text documents

to <u>predictions</u>

- numerical value (probability stock price increases)
- label (is the image a cat? does the image contain a car?)
- decision (turn car left, rotate robotic arm)

## MACHINE LEARNING MODEL

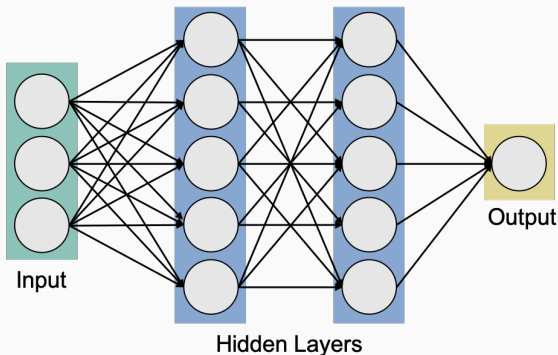Let $M_\mathbf{x}$ be a model with parameters $\mathbf{x} = \{x_1, \ldots, x_k\}$, which takes as input a data vector $\mathbf{a}$ and outputs a prediction.

**Example:**

$$M_\mathbf{x}(\mathbf{a}) = \text{sign}(\mathbf{a}^T \mathbf{x})$$

**Example:**



Input

Hidden Layers

Output

$\mathbf{x} \in \mathbb{R}^{(\# \text{ of connections})}$ is the parameter vector containing all the network weights.

## SUPERVISED LEARNING

Classic approach in <u>supervised learning</u>: Find a model that works well on data that you already have the answer for (labels, values, classes, etc.).

- Model $M_{\mathbf{x}}$ parameterized by a vector of numbers $\mathbf{x}$.
- Dataset $\mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(n)}$ with outputs $y^{(1)}, \ldots, y^{(n)}$.

    Want to find $\hat{\mathbf{x}}$ so that $M_{\hat{\mathbf{x}}}(\mathbf{a}^{(i)}) \approx y^{(i)}$ for $i \in 1, \ldots, n$.

**How do we turn this into a function minimization problem?**

## LOSS FUNCTION

**Loss function** $L(M_x(a), y)$: Some measure of distance between prediction $M_x(a)$ and target output $y$. Increases if they are further apart.

- Squared ($\ell_2$) loss: $|M_x(a) - y|^2$
- Absolute deviation ($\ell_1$) loss: $|M_x(a) - y|$
- Hinge loss: $1 - y \cdot M_x(a)$
- Cross-entropy loss (log loss).
- Etc.

## EMPIRICAL RISK MINIMIZATION

**Empirical risk minimization**:

$$f(\mathbf{x}) = \sum_{i=1}^{n} L\left(M_{\mathbf{x}}(\mathbf{a}^{(i)}), y^{(i)}\right)$$

Solve the optimization problem $\min_{\mathbf{x}} f(\mathbf{x})$.

## EXAMPLE: LINEAR REGRESSION

- $M_{\mathbf{x}}(\mathbf{a}) = \mathbf{x}^T \mathbf{a}$. $\mathbf{x}$ contains the regression coefficients.
- $L(z, y) = |z - y|^2$.
- $f(\mathbf{x}) = \sum_{i=1}^{n} |\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}|^2$

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

where $\mathbf{A}$ is a matrix with $\mathbf{a}^{(i)}$ as its $i^{\text{th}}$ row and $\mathbf{y}$ is a vector with $y^{(i)}$ as its $i^{\text{th}}$ entry.

$$\min_x \; |Ax - y|_2^2 + \alpha|x|_1^2$$

The choice of algorithm to minimize $f(\mathbf{x})$ will depend on:
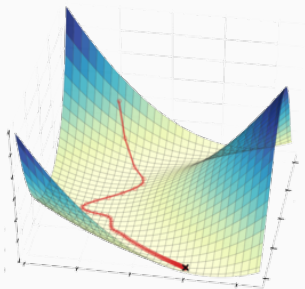
- The form of $f(\mathbf{x})$ (is it linear, is it quadratic, does it have finite sum structure, etc.)

- If there are any additional constraints imposed on $\mathbf{x}$. E.g. $\|\mathbf{x}\|_2 \leq c$.

**What are some example algorithms for continuous optimization?**

LP $\Longleftrightarrow$ simplex

Ellipoids

semi-definite programs

interior point methods

# GRADIENT DESCENT

**Gradient descent:** A greedy algorithm for minimizing functions of multiple variables that often works amazingly well.



(and sometimes we can prove it works)

## CALCULUS REVIEW

For $i = 1, \ldots, d$, let $x_i$ be the $i^{\text{th}}$ entry of $\mathbf{x}$. Let $\mathbf{e}^{(i)}$ be the $i^{\text{th}}$ standard basis vector.

**Partial derivative:**

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{e}^{(i)}) - f(\mathbf{x})}{t}$$

**Directional derivative:**

$$D_{\mathbf{v}} f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}$$

**Gradient**:

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

**Directional derivative:**

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

## FIRST ORDER OPTIMIZATION

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 $\mathbf{x}$.

We view the implementation of these oracles as black-boxes, but they can often require a fair bit of computation.

$n \gg d$

$A \in \mathbb{R}^{n \times d}$

**Linear least-squares regression**:

- Given $\mathbf{a}^{(1)}, \ldots \mathbf{a}^{(n)} \in \mathbb{R}^d$, $y^{(1)}, \ldots y^{(n)} \in \mathbb{R}$.
- Want to minimize:

$(Ax - y)^T (Ax - y)$

$$f(\mathbf{x}) = \sum_{i=1}^{n} \left( \mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)} \right)^2 = \|\mathbf{Ax} - \mathbf{y}\|_2^2.$$

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^{n} 2 \left( \mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)} \right) \cdot a_j^{(i)} = (2\mathbf{Ax} - \mathbf{y})^T \boldsymbol{\alpha}^{(j)}$$

$2 A^T A x - 2 A^T y$

where $\boldsymbol{\alpha}^{(j)}$ is the $j^{\text{th}}$ column of $\mathbf{A}$.

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{y})$$

**What is the time complexity of a gradient oracle for $\nabla f(\mathbf{x})$?**

$A^T A$ in time $O(\min(nd^2), (n^2 d))$

$$\nabla f(x) = [100, -1]$$
$$\text{could set } v = (\tfrac{1}{v_1}, \tfrac{1}{v_2})$$
$$v = (-1, 0)$$

**Greedy approach:** Given a starting point $\mathbf{x}$, make a small adjustment that decreases $f(\mathbf{x})$. In particular, $\mathbf{x} \leftarrow \mathbf{x} + \eta\mathbf{v}$ and $f(\mathbf{x}) \leftarrow f(\mathbf{x} + \eta\mathbf{v})$.

What property do I want in **v**?

**Leading question:** When $\eta$ is small, what's an approximation for $f(\mathbf{x} + \eta\mathbf{v}) - f(\mathbf{x})$?

$$f(\mathbf{x} + \eta\mathbf{v}) - f(\mathbf{x}) \approx \eta \, \nabla f(x)^T \cdot v$$

$$V = -\eta f(x)^T$$

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

So:

$$f(\mathbf{x} + \eta\mathbf{v}) - f(\mathbf{x}) \approx \eta \nabla f(x)^T v$$

**How should we choose v so that $f(\mathbf{x} + \eta\mathbf{v}) < f(\mathbf{x})$?**

$$v = -\frac{\nabla f(x)^T}{|\nabla f(x)^T|_2}$$
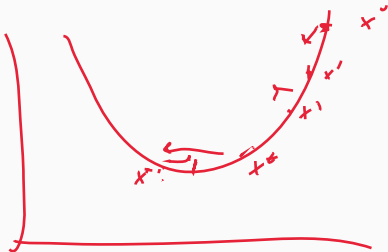
## GRADIENT DESCENT

**Prototype algorithm:**

- Choose starting point $\mathbf{x}^{(0)}$.
- For $i = 0, \ldots, T$:
    - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\mathbf{x}^{(T)}$.

$\eta$ is a step-size parameter, which is often adapted on the go. For now, assume it is fixed ahead of time.
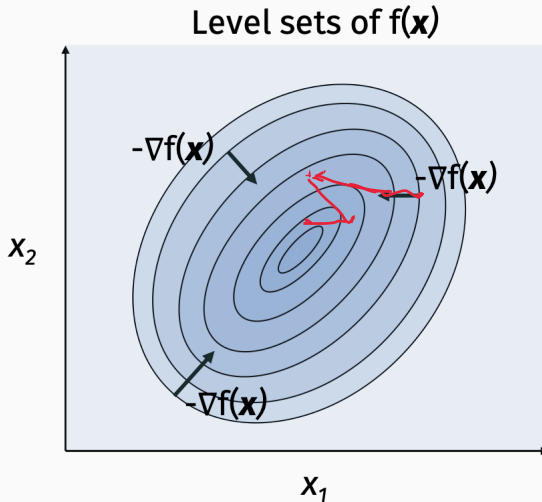
**1 dimensional example:**

**2 dimensional example:**



Level sets of f($\boldsymbol{x}$)

**For a convex function** $f(\mathbf{x})$**:** For sufficiently small $\eta$ and a sufficiently large number of iterations $T$, gradient descent will converge to a **near global minimum**:

$$f(\mathbf{x}^{(T)}) \leq f(\mathbf{x}^*) + \epsilon.$$

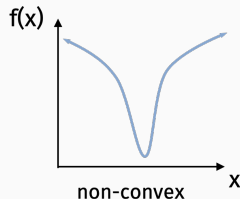Examples: least squares regression, logistic regression, kernel regression, SVMs.

**For a non-convex function** $f(\mathbf{x})$**:** For sufficiently small $\eta$ and a sufficiently large number of iterations $T$, gradient descent will converge to a **near stationary point**:

$$\|\nabla f(\mathbf{x}^{(T)})\|_2 \leq \epsilon.$$

Examples: neural networks, matrix completion problems, mixture models.

One issue with non-convex functions is that they can have **local minima**. Even when they don't, convergence analysis requires different assumptions than convex functions.

## APPROACH FOR THIS UNIT

We care about how fast gradient descent and related methods converge, not just that they do converge.

- Bounding iteration complexity requires placing some assumptions on $f(\mathbf{x})$.
- Stronger assumptions lead to better bounds on the convergence.

Understanding these assumptions can help us design faster variants of gradient descent (there are many!).

Today, we will start with **convex functions** only.

# CONVEXITY

## Definition (Convex)

A function $f$ is convex iff for any $\mathbf{x}, \mathbf{y}, \lambda \in [0, 1]$:

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



$h = (1 - \lambda)x + \lambda y$

$\mathbf{h} = (1-\lambda)\mathbf{x} + \lambda\mathbf{y}$

f(**y**)

f(**h**)

f(**x**)

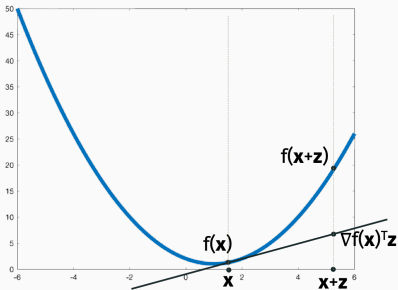**x**    **h**    **y**

**Definition (Convex)**

A function $f$ is convex if and only if for any $\mathbf{x}, \mathbf{y}$:

$$f(\mathbf{x} + \mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{z}$$

$z = x - y$

Equivalently:

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{y})$$

$f(y) = \sqrt{y}$

# GRADIENT DESCENT ANALYSIS

**Assume:** $\left| f(x) - f(y) \right|_2 \leq G \cdot |x - y|_2$

- $f$ 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 $i = 0, \ldots, T$:
    - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg\min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)})$.

**Claim (GD Convergence Bound)**

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



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

## "FUNDAMENTAL THEOREM OF OPTIMIZATION"

**Fact:** For any two vectors $v, u$ of the same dimension, we have:

$$v^T u = \langle v, u \rangle = \frac{1}{2} \left( \|v\|_2^2 + \|u\|_2^2 - \|u - v\|_2^2 \right)$$

**Proof:** Recall $\|u - v\|_2^2 = \|u\|_2^2 + \|v\|_2^2 - 2\langle u, v \rangle$

Inner products can be written as a sum of norms!

**Claim (GD Convergence Bound)**

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

$$x^{(i+1)} \leftarrow x^{(i)} - \eta \nabla f(x^{(i)})$$

**Proof:** For all $i = 0, \ldots, T$:

$$f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x}^{(i)} - \mathbf{x}^*) \quad \text{convexity}$$

$$= \frac{1}{\eta} \langle \mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}, \mathbf{x}^{(i)} - \mathbf{x}^* \rangle \quad \begin{array}{l} \text{gradient} \\ \text{update} \end{array}$$

## GRADIENT DESCENT ANALYSIS

**Claim (GD Convergence Bound)**

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

**Proof:** For all $i = 0, \ldots, T$:

$$
\begin{aligned}
f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x}^{(i)} - \mathbf{x}^*) \\
&= \frac{1}{\eta} \langle \mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}, \, \mathbf{x}^{(i)} - \mathbf{x}^* \rangle \\
&\leq \frac{1}{2\eta} (\|\mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}\|_2^2 + \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2)
\end{aligned}
$$

**Claim (GD Convergence Bound)**

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

**Proof:** For all $i = 0, \ldots, T$:

$$
\begin{aligned}
f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x}^{(i)} - \mathbf{x}^*) \\
&= \frac{1}{\eta} \langle \mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}, \mathbf{x}^{(i)} - \mathbf{x}^* \rangle \\
&\leq \frac{1}{2\eta} (\|\mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}\|_2^2 + \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2) \\
&\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{1}{2\eta} \|\eta \nabla f(\mathbf{x}^{(i)})\|_2^2
\end{aligned}
$$

**Claim (GD Convergence Bound)**

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

**Proof:** For all $i = 0, \ldots, T$:

$$\|\nabla f(\mathbf{x})\|_2 \leq G$$

$$
\begin{aligned}
f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x}^{(i)} - \mathbf{x}^*) \\
&= \frac{1}{\eta} \langle \mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}, \ \mathbf{x}^{(i)} - \mathbf{x}^* \rangle \\
&\leq \frac{1}{2\eta} (\|\mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}\|_2^2 + \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2) \\
&\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{1}{2\eta} \|\underbrace{\eta \nabla f(\mathbf{x}^{(i)})}_{\leq \eta \, G}\|_2^2 \\
&\leq \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}
\end{aligned}
$$

34

## Claim (GD Convergence Bound)

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

**Proof:** For all $i = 0, \ldots, T$,

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

**Telescoping sum:**

$$\sum_{i=0}^{T-1} \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \frac{\overbrace{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}^{< R} - \overbrace{\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}^{< 0}}{2\eta} + \frac{T \eta G^2}{2}$$

$$\frac{1}{T} \sum_{i=0}^{T-1} \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \underbrace{\frac{R^2}{2T\eta}}_{< \frac{\epsilon}{2}} + \underbrace{\frac{\eta G^2}{2}}_{< \frac{\epsilon}{2}}$$

$$\eta G \leq \frac{R}{\sqrt{T}} G$$
$$< \epsilon$$

35

## GRADIENT DESCENT ANALYSIS

**Claim (GD Convergence Bound)**

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

**Final step:**

$$\frac{1}{T} \sum_{i=0}^{T-1} \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \epsilon$$

$$\left[ \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)}) \right] - f(\mathbf{x}^*) \leq \epsilon$$

We always have that $\min_i f(\mathbf{x}^{(i)}) \leq \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)})$, so this is what we return:

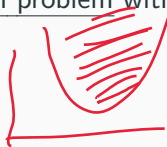$$f(\hat{\mathbf{x}}) = \min_{i \in 1, \dots, T} f(\mathbf{x}^{(i)}) \leq f(\mathbf{x}^*) + \epsilon.$$

$S = \{x \mid |x|_2 < R\}$

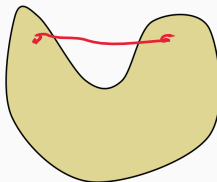**Typical goal**: Solve a <u>convex minimization problem with</u> additional <u>convex constraints</u>.

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

where $\mathcal{S}$ is a **convex set**.

Which of these is convex?

# CONSTRAINED CONVEX OPTIMIZATION



**Definition (Convex set)**

A set $\mathcal{S}$ is convex if for any $\mathbf{x}, \mathbf{y} \in \mathcal{S}, \lambda \in [0, 1]$:

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

**Gradient descent:**

- For $i = 0, \ldots, T$:
    - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$
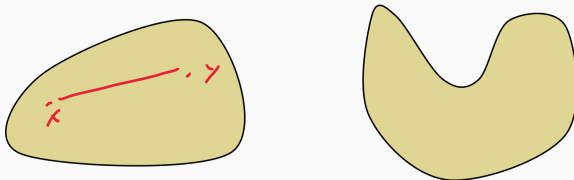- Return $\hat{\mathbf{x}} = \arg \min_i f(\mathbf{x}^{(i)})$.

Even if we start with $\mathbf{x}^{(0)} \in \mathcal{S}$, there is no guarantee that $\mathbf{x}^{(0)} - \eta \nabla f(\mathbf{x}^{(0)})$ will remain in our set.

**Extremely simple modification:** Force $\mathbf{x}^{(i)}$ to be in $\mathcal{S}$ by **projecting** onto the set.

## CONSTRAINED FIRST ORDER OPTIMIZATION

Given a function $f$ to minimize and a convex constraint set $\mathcal{S}$, assume we have:

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

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

$$P_S(x) = \begin{cases} x/|x|_2 & \text{if } x \notin S \\ x & \text{o.w.} \end{cases}$$

- How would you implement $P_S$ for $S = \{\mathbf{y} : \|\mathbf{y}\|_2 \leq 1\}$.

- How would you implement $P_S$ for $S = \{\mathbf{y} : \mathbf{y} = \mathbf{Q}\mathbf{z}\}$.



$x$    $a z$

$\mathbf{x}$

$P_S(\mathbf{x})$

$$\min_{z \in \mathbb{R}^k} |Qz - x|_2$$

41

## PROJECTED GRADIENT DESCENT

Given function $f(\mathbf{x})$ and set $\mathcal{S}$, such that $\|\nabla f(\mathbf{x})\|_2 \leq G$ for all $\mathbf{x} \in \mathcal{S}$ and starting point $\mathbf{x}^{(0)}$ with $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq R$.

**Projected gradient descent:**

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

**Claim (PGD Convergence Bound)**

If $f, \mathcal{S}$ are convex and $T \geq \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$.
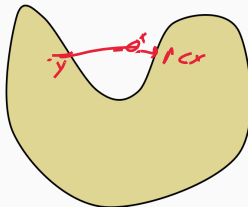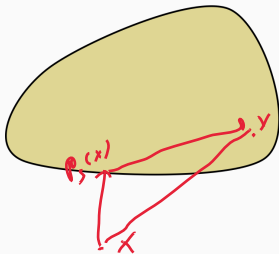
# PROJECTED GRADIENT DESCENT ANALYSIS

Analysis is almost identical to standard gradient descent! We just need one additional claim:

**Claim (Contraction Property of Convex Projection)**

If $\mathcal{S}$ is convex, then for <u>any</u> $\mathbf{y} \in \mathcal{S}$,

$$\|\mathbf{y} - P_{\mathcal{S}}(\mathbf{x})\|_2 \leq \|\mathbf{y} - \mathbf{x}\|_2.$$

**Claim (PGD Convergence Bound)**

If $f, \mathcal{S}$ are convex and $T \geq \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \epsilon$.

$\mathbf{x}^{(i+1)} = P_{\mathcal{S}}(\mathbf{x})$

**Claim 1:** For all $i = 0, \ldots, T$,

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

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

**Same telescoping sum argument:**

$$\left[ \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)}) \right] - f(\mathbf{x}^*) \leq \frac{R^2}{2T\eta} + \frac{\eta G^2}{2}.$$

**Conditions:**

- **Convexity:** $f$ is a convex function, $\mathcal{S}$ is a convex set.
- **Bounded initial distant:**

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

- **Bounded gradients (Lipschitz function)**:

$$\|\nabla f(\mathbf{x})\|_2 \leq G \text{ for all } \mathbf{x} \in \mathcal{S}.$$

**Theorem**

*GD Convergence Bound] (Projected) Gradient Descent returns $\hat{\mathbf{x}}$ with $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}) + \epsilon$ after*

$$T = \frac{R^2 G^2}{\epsilon^2} \text{ iterations.}$$

**break**

$$L(x, y) = x y^T$$
$$\nabla_x L = y^T$$
$$\|\nabla_x L\|_2 \leq \|y\|_2$$

Can our convergence bound be tightened for certain functions?
Can it guide us towards faster algorithms?

**Goals:**

- Improve $\epsilon$ dependence below $1/\epsilon^2$.
  - Ideally $1/\epsilon$ or $\log(1/\epsilon)$.
- Reduce or eliminate dependence on $G$ and $R$.

**Definition ($\beta$-smoothness)**

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

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

$|f(x) - f(y)| \leq G|x - y|$

After some calculus (see Lem. 3.4 in **Bubeck's book**), this implies:

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



For a scalar valued function $f$, equivalent to $f''(x) \leq \beta$.

Recall from definition of convexity that:

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$

So now we have an upper and lower bound.

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

Previously learning rate/step size $\eta$ depended on $G$. Now choose it based on $\beta$:

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

**Progress per step of gradient descent:**

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

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

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

49

## CONVERGENCE GUARANTEE

**Theorem (GD convergence for $\beta$-smooth functions.)**

*Let $f$ be a $\beta$ 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 $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon$.

Complete proof in Theorem 3.5 of **Bubeck's book**

**Definition ($\alpha$-strongly convex)**

$< \frac{\mu}{2} |x - y|_2^2$

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

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

$\alpha$ is a parameter that will depend on our function.

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

## GD FOR STRONGLY CONVEX FUNCTION

**Gradient descent for strongly convex functions:**

- Choose number of steps $T$.

- For $i = 1, \ldots, T$:
    - $\eta = \frac{2}{\alpha \cdot (i+1)}$
    - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$

- Return $\hat{\mathbf{x}} = \arg\min_{\mathbf{x}^{(i)}} f(\mathbf{x}^{(i)})$.

**Theorem (GD convergence for $\alpha$-strongly convex functions.)**

*Let $f$ be an $\alpha$-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:*

$$f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{2G^2}{\alpha(T-1)}$$

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

## CONVERGENCE GUARANTEE

What if $f$ is both $\beta$-smooth and $\alpha$-strongly convex?

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

## CONVERGENCE GUARANTEE

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

**Theorem (GD for $\beta$-smooth, $\alpha$-strongly convex.)**

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

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 \leq e^{-(T-1)\frac{\alpha}{\beta}}\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2$$

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

**Is it better if $\kappa$ is large or small?**

**Converting to more familiar form:** Using that fact the $\nabla f(\mathbf{x}^*) = \mathbf{0}$ along with

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

we have:

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2^2 \leq \frac{2}{\alpha}\left[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)\right]$$

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

$$f(x^T) - f(\vec{x^*}) < \frac{\ell}{2}|x^T - \vec{x^*}|$$

$$< \frac{\ell}{2} e^{-T\frac{\alpha}{\ell}}|x' - x^+|$$

$$< \frac{\ell}{2\alpha} e^{-T\frac{\alpha}{\ell}}|f(x') - f(x)|$$

**Corollary (GD for $\beta$-smooth, $\alpha$-strongly convex.)**

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \frac{\beta}{\alpha} e^{-(T-1)\frac{\alpha}{\beta}} \cdot \left[ f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$

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

$$f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \epsilon \left[ f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$

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

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

Let $f$ be a twice differentiable function from $\mathbb{R}^d \to \mathbb{R}$. Let the **Hessian** $\mathbf{H} = \nabla^2 f(\mathbf{x})$ contain all of its second derivatives at a point $\mathbf{x}$. So $\mathbf{H} \in \mathbb{R}^{d \times d}$. We have:

$$\mathbf{H}_{i,j} = \left[\nabla^2 f(\mathbf{x})\right]_{i,j} = \frac{\partial^2 f}{\partial x_i x_j}.$$
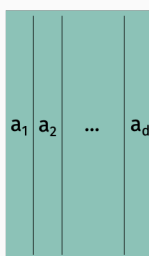
For vector $\mathbf{x}, \mathbf{v}$:

$$\nabla f(\mathbf{x} + t\mathbf{v}) \approx \nabla f(\mathbf{x}) + t \left[\nabla^2 f(\mathbf{x})\right] \mathbf{v}.$$

Let $f$ be a twice differentiable function from $\mathbb{R}^d \to \mathbb{R}$. Let the **Hessian** $\mathbf{H} = \nabla^2 f(\mathbf{x})$ contain all of its second derivatives at a point $\mathbf{x}$. So $\mathbf{H} \in \mathbb{R}^{d \times d}$. We have:

$$\mathbf{H}_{i,j} = \left[\nabla^2 f(\mathbf{x})\right]_{i,j} = \frac{\partial^2 f}{\partial x_i x_j}.$$

**Example:** Let $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$. Recall that $\nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$.

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



$a_1$ $a_2$ ... $a_d$

X

A

b

## HESSIAN MATRICES AND POSITIVE SEMIDEFINITENESS

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

**Definition (Positive Semidefinite (PSD))**

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is <u>positive semidefinite</u> (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 <u>partial ordering</u> on matrices.

## HESSIAN MATRICES AND POSITIVE SEMIDEFINITENESS

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

**Definition (Positive Semidefinite (PSD))**

A square, symmetric matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$ is <u>positive semidefinite</u> (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: $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, $\mathbf{H} = \nabla^2 f(\mathbf{x}) = 2\mathbf{A}^T \mathbf{A}$ for all $\mathbf{x}$. Is $\mathbf{H}$ PSD?

## THE LINEAR ALGEBRA OF CONDITIONING

If $f$ is $\beta$-smooth and $\alpha$-strongly convex then at any point $\mathbf{x}$, $\mathbf{H} = \nabla^2 f(\mathbf{x})$ satisfies:

$$\alpha \mathbf{I}_{d \times d} \preceq \mathbf{H} \preceq \beta \mathbf{I}_{d \times d},$$

where $\mathbf{I}_{d \times d}$ 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

$$\alpha \mathbf{I}_{d \times d} \preceq \mathbf{H} \preceq \beta \mathbf{I}_{d \times d}.$$

Equivalently for any $\mathbf{z}$,

$$\alpha \|\mathbf{z}\|_2^2 \leq \mathbf{z}^T \mathbf{H} \mathbf{z} \leq \beta \|\mathbf{z}\|_2^2.$$
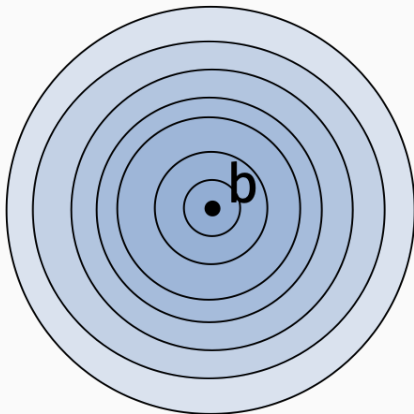
## SIMPLE EXAMPLE

Let $f(\mathbf{x}) = \|\mathbf{Dx} - \mathbf{b}\|_2^2$ where $\mathbf{D}$ is a diagaonl 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 $\alpha, \beta$ for this problem?**
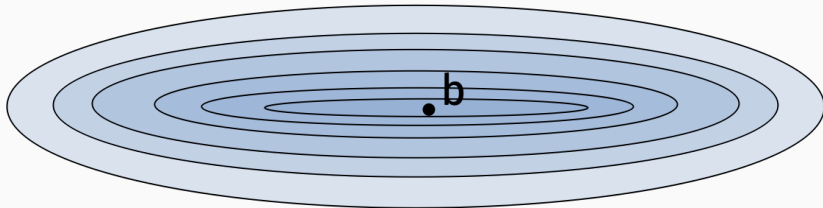
$$\alpha \|\mathbf{z}\|_2^2 \leq \mathbf{z}^T \mathbf{H} \mathbf{z} \leq \beta \|\mathbf{z}\|_2^2$$

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

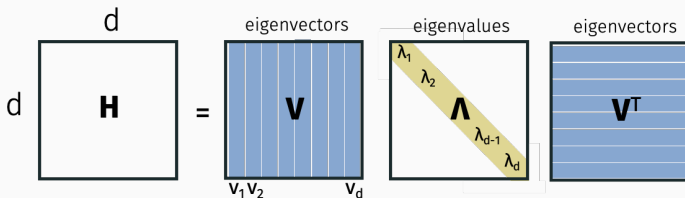Level sets of $\|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$ when $d_1^2 = \frac{1}{3}, d_2^2 = 2$.

## EIGENDECOMPOSITION VIEW

Any symmetric matrix $\mathbf{H}$ has an <u>orthogonal</u>, real valued eigendecomposition.
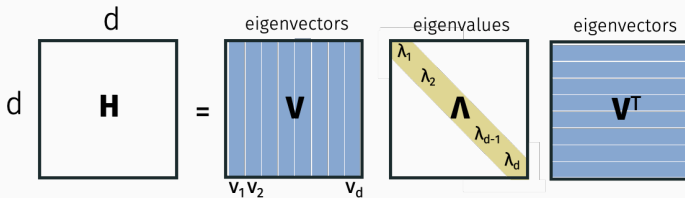


Here $\mathbf{V}$ is square and orthogonal, so $\mathbf{V}^T\mathbf{V} = \mathbf{V}\mathbf{V}^T = \mathbf{I}$. And for each $\mathbf{v}_i$, we have:

$$\mathbf{H}\mathbf{v}_i = \lambda_i\mathbf{v}_i.$$

By definition, that's what makes $\mathbf{v}_1, \ldots, \mathbf{v}_d$ eigenvectors.
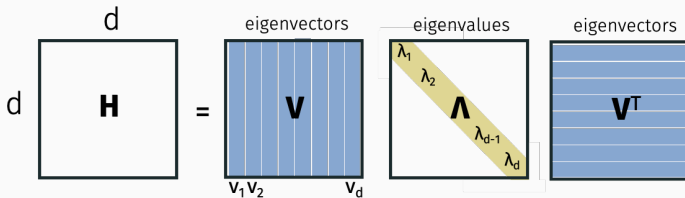
## EIGENDECOMPOSITION VIEW

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



**Claim:** $\mathbf{H}$ is PSD $\Leftrightarrow \lambda_1, ..., \lambda_d \geq 0$.
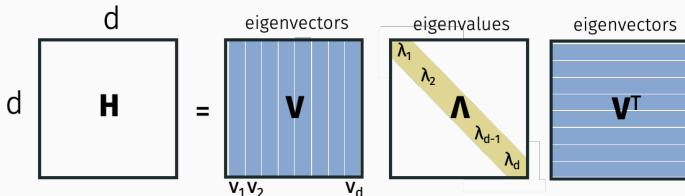
## EIGENDECOMPOSITION VIEW

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



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

Recall $VV^T = V^TV = I$.



In other words, if we let $\lambda_{\max}(H)$ and $\lambda_{\min}(H)$ be the smallest and largest eigenvalues of $H$, then for all $z$ we have:

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

## EIGENDECOMPOSITION VIEW

If the maximum eigenvalue of $\mathbf{H} = \nabla^2 f(\mathbf{x}) = \beta$ and the minimum eigenvalue of $\mathbf{H} = \nabla^2 f(\mathbf{x}) = \alpha$ then $f(\mathbf{x})$ is $\beta$-smooth and $\alpha$-strongly convex.

$$\lambda_{\max}(\mathbf{H}) = \beta$$
$$\lambda_{\min}(\mathbf{H}) = \alpha$$

## POLYNOMIAL VIEW POINT

**Theorem (GD for $\beta$-smooth, $\alpha$-strongly convex.)**

*Let $f$ be a $\beta$-smooth and $\alpha$-strongly convex function. If we run GD for $T$ steps (with step size $\eta = \frac{2}{\beta}$) we have:*

$$\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2 \leq e^{-T/\kappa}\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2$$

**Goal: Prove for $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$.**

Let $\lambda_{\max} = \lambda_{\max}(\mathbf{A}^T\mathbf{A})$. Gradient descent update is:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{2\,\lambda_{\max}}2\mathbf{A}^T(\mathbf{A}\mathbf{x}^{(t)} - \mathbf{b})$$

**Richardson Iteration view:**

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

What is the maximum eigenvalue of the symmetric matrix $\left(\mathbf{I} - \frac{1}{\lambda_{\max}}\mathbf{A}^T\mathbf{A}\right)$ in terms of the eigenvalues $\lambda_{\max} = \lambda_1 \geq \ldots \geq \lambda_d = \lambda_{\min}$ of $\mathbf{A}^T\mathbf{A}$?

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

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

So we have $\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2 \leq$

# IMPROVING GRADIENT DESCENT

We now have a pretty good understanding of gradient descent.

**Number of iterations for $\epsilon$ error:**

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