Gradient Descent

Yingzhen Li

Department of Computing Imperial College London

October 26, 2021

✓@liyzhen2 yingzhen.li@imperial.ac.uk

Reading for this week

Read MML book: Sections 4.1 - 4.4, 7.1, Chapter 9 up to 9.2.3 Do MML book exercises: Exercises 4.1 - 4.7 An extra exercise will be uploaded to course materials.

Optimisation Problems



Remember our problem:

- · Find a curve that predicts well even for unseen inputs
- Start by minimising loss on training points:

$$L(\boldsymbol{\theta}) = \sum_{n} (f(\boldsymbol{x}_{n}; \boldsymbol{\theta}) - y_{n})^{2}$$
(1)

Formulating optimisation problems

- Define **objective function** $L : \mathbb{R}^D \to \mathbb{R}$
- Unconstrained minimisation, state which variable you want to optimise over

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \tag{2}$$

How do we solve optimisation? We need an algorithm.

How do we solve optimisation? We need an **algorithm**. Gradient-based optimisation is a **class** of methods which all

Pick some starting point θ₀.

- Pick some starting point θ₀.
- Iteratively update the parameters, resulting in a sequence of solutions θ₁,..., θ_T.

- Pick some starting point θ₀.
- Iteratively update the parameters, resulting in a sequence of solutions θ₁,..., θ_T.
- Choose the update of the parameter by computing the **gradient** $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$.

- Pick some starting point θ₀.
- Iteratively update the parameters, resulting in a sequence of solutions θ₁,..., θ_T.
- Choose the update of the parameter by computing the **gradient** $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$.
- Usually a **stopping criterion** (e.g. iteration budget, time budget, gradient size, ...)

Gradient descent

Algorithm: Gradient Descent

Define starting point θ_0 , sequence of step sizes γ_t , set $t \leftarrow 0$.

1. Set
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t), t \leftarrow t+1$$

2. Repeat 1 until stopping criterion.



Gradient Descent

Fitting linear regression models:

• Dataset:
$$\mathcal{D} = \{\mathbf{X}, \mathbf{y}\},\$$

 $\mathbf{X} = [\mathbf{x}_1, ..., \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D},\$
 $\mathbf{y} = [y_1, ..., y_N]^\top \in \mathbb{R}^{N \times 1}$

• Goal: find $\boldsymbol{\theta} \in \mathbb{R}^{D \times 1}$ such that

 $y\approx X\pmb{\theta}$



A typical linear regression model:

- $x \in \mathbb{R}^{D \times 1}$: input features; $y \in \mathbb{R}$: output value
- Model and loss:

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \boldsymbol{x}^{\top} \boldsymbol{\theta}, \quad \boldsymbol{y} = f(\boldsymbol{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon}, \ \boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2)$$
$$L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \sum_n (f(\boldsymbol{x}_n, \boldsymbol{\theta}) - \boldsymbol{y}_n)^2$$

A typical linear regression model:

- $x \in \mathbb{R}^{D \times 1}$: input features; $y \in \mathbb{R}$: output value
- Model and loss:

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \boldsymbol{x}^{\top} \boldsymbol{\theta}, \quad \boldsymbol{y} = f(\boldsymbol{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon}, \ \boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2)$$
$$L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \sum_n (f(\boldsymbol{x}_n, \boldsymbol{\theta}) - y_n)^2$$

• Rewriting the loss in matrix form:

$$L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2$$

Optimal solution of θ :

$$\boldsymbol{\theta}^* = \operatorname*{arg\,min}_{\boldsymbol{\theta}} \quad L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2$$

• Gradient of the loss $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$:

Optimal solution of θ :

$$\boldsymbol{\theta}^* = \operatorname*{arg\,min}_{\boldsymbol{\theta}} \quad L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2$$

• Gradient of the loss $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$:

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

• Setting $\nabla_{\theta} L(\theta) = 0$:

Optimal solution of θ :

$$\boldsymbol{\theta}^* = \operatorname*{arg\,min}_{\boldsymbol{\theta}} \quad L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_2^2$$

• Gradient of the loss $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$:

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

• Setting $\nabla_{\theta} L(\theta) = 0$:

$$\Rightarrow \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta}^* = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y}$$
$$\Rightarrow \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Gradient descent to find θ^* :

Assume constant step-sizes $\gamma_t = \gamma$:

1. Define **starting point** θ_0 , set $t \leftarrow 0$

2. Set
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t), t \leftarrow t+1$$

$$egin{aligned} oldsymbol{ heta}_{t+1} &= oldsymbol{ heta}_t - \gamma_t
abla_{oldsymbol{ heta}} L(oldsymbol{ heta}_t) \ &= oldsymbol{ heta}_t - \gamma rac{1}{\sigma^2} \mathbf{X}^{ op} (\mathbf{X} oldsymbol{ heta}_t - \mathbf{y}) \end{aligned}$$

Gradient descent to find θ^* :

Assume constant step-sizes $\gamma_t = \gamma$:

1. Define starting point θ_0 , set $t \leftarrow 0$

2. Set
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t), t \leftarrow t+1$$

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t - \gamma \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta}_t - \mathbf{y}) \\ &= (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

3. Repeat 1 until stopping criterion.

Gradient descent to find θ^* :

Assume constant step-sizes $\gamma_t = \gamma$:

• GD returns the following iterative updates:

$$\boldsymbol{\theta}_{t+1} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{y}$$

• Solving this iterative update returns:

Gradient descent to find θ^* :

Assume constant step-sizes $\gamma_t = \gamma$:

• GD returns the following iterative updates:

$$\boldsymbol{\theta}_{t+1} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{y}$$

Solving this iterative update returns:

$$\boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*, \quad \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

• GD converges $(\boldsymbol{\theta}_t \to \boldsymbol{\theta}^*)$ if $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \to \mathbf{0}$

11











 $Qz = x \qquad \Rightarrow \qquad \begin{array}{l} \text{Find vector } x \text{ such that} \\ z = Q^{-1}x \\ Q = \begin{pmatrix} q_{11} & q_{12} \\ q_{21} & q_{22} \end{pmatrix} = (q_1, q_2) \qquad Q^{-1} = \begin{pmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{pmatrix} = \begin{pmatrix} q_1^{\mathsf{T}} \\ q_2^{\mathsf{T}} \end{pmatrix} \\ \text{Left multiply } Q \text{ on } z: \\ \text{Change basis from } \{q_1, q_2\} \text{ to } \{e_1, e_2\} \end{array}$









$$\begin{aligned} \mathbf{x}' &= A\mathbf{x}, A = Q\Lambda Q^{-1} \implies \mathbf{x}' = \mathbf{z}_1' q_1 + \mathbf{z}_2' q_2 \\ Q^{-1} &= \begin{pmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{pmatrix} = \begin{pmatrix} q_1^{\mathsf{T}} \\ q_2^{\mathsf{T}} \end{pmatrix} \qquad \mathbf{z}_1' = \lambda_1 q_1^{\mathsf{T}} \mathbf{x} \\ \mathbf{z}_2' &= \lambda_2 q_2^{\mathsf{T}} \mathbf{x} \\ \Lambda &= \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \qquad \text{(for orthonormal basis, } Q^{-1} = Q^{\mathsf{T}} \end{aligned}$$

Consider x' = Ax for $A \in \mathbb{R}^{D \times D}$:

- Can we find $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$ (and how)?
- Key idea: find x such that
 x' = Ax and x align on the same line.

i.e.
$$\mathbf{x}' = \lambda \mathbf{x}$$
 for some $\lambda \neq 0$





$$\begin{aligned} \mathbf{x}' &= A\mathbf{x}, A = Q\Lambda Q^{-1} \implies \mathbf{x}' = \mathbf{z}_1' q_1 + \mathbf{z}_2' q_2 \\ Q^{-1} &= \begin{pmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{pmatrix} = \begin{pmatrix} q_1^{\mathsf{T}} \\ q_2^{\mathsf{T}} \end{pmatrix} \qquad \mathbf{z}_1' = \lambda_1 q_1^{\mathsf{T}} \mathbf{x} \\ \mathbf{z}_2' &= \lambda_2 q_2^{\mathsf{T}} \mathbf{x} \\ \Lambda &= \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \qquad \text{(for orthonormal basis, } Q^{-1} = Q^{\mathsf{T}} \text{)} \end{aligned}$$

For i = 1, ..., D, if $x = q_i$, then $x' = \mathbf{A}x = \lambda_i q_i$.

Eigenvalue decomposition for a matrix $\mathbf{A} \in \mathbb{R}^{D \times D}$: Find scalar λ and vector \boldsymbol{q} such that

$$\mathbf{A}\boldsymbol{q}=\lambda\boldsymbol{q}$$

If solutions exist: $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$, $\Lambda = \text{diag}(\lambda_1, ..., \lambda_D)$

Eigenvalue decomposition for a matrix $\mathbf{A} \in \mathbb{R}^{D \times D}$: Find scalar λ and vector \boldsymbol{q} such that

$$\mathbf{A}\boldsymbol{q}=\lambda\boldsymbol{q}$$

If solutions exist: $\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}$, $\Lambda = \text{diag}(\lambda_1, ..., \lambda_D)$

If **A** is symmetric:

- there are *D* pairs of solutions $(\lambda_d, \boldsymbol{q}_d)$ such that
 - Column vectors in Q = (q₁, ..., q_D) form an orthonormal basis (so Q⁻¹ = Q^T)
 - $\bullet \ \lambda_1 \geqslant \lambda_2 \geqslant \ldots \geqslant \lambda_D$
 - If **A** is positive semi-definite: $\lambda_D \ge 0$

15

How to find $\mathbf{A} = \mathbf{Q} \wedge \mathbf{Q}^{-1}$:

- $\mathbf{A}q = \lambda q \quad \Rightarrow \quad (\mathbf{A} \lambda \mathbf{I})q = \mathbf{0}$
- Assume $q \neq 0$: find λ such that

$$det(\mathbf{A} - \lambda \mathbf{I}) = 0$$

• Once λ is solved, plugging it back and solve for q

16



Why solving $det(\mathbf{A} - \lambda \mathbf{I})$ w.r.t. λ :

► |*det*(**A**)| tells how volume is scaled by linear transform
Eigen decomposition



Why solving $det(\mathbf{A} - \lambda \mathbf{I})$ w.r.t. λ :

- |*det*(**A**)| tells how volume is scaled by linear transform
- If *det*(**A** − λ**I**) = 0: some subspace in ℝ^D is squashed to {**0**}
 ⇒ there exists *q* ≠ **0** such that (**A** − λ**I**)*q* = **0**

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$:

• If λ is an eigenvalue of **A**, then λ^t is an eigenvalue of **A**^t:

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$:

• If λ is an eigenvalue of **A**, then λ^t is an eigenvalue of **A**^t:

$$\mathbf{A}^{t} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}\mathbf{Q}\Lambda\mathbf{Q}^{-1}\cdots\mathbf{Q}^{-1}\Lambda\mathbf{Q}^{-1} = \mathbf{Q}\Lambda^{t}\mathbf{Q}^{-1},$$
$$\Lambda^{t} = \operatorname{diag}(\lambda_{1}^{t},...,\lambda_{D}^{t})$$

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$:

• If λ is an eigenvalue of **A**, then λ^t is an eigenvalue of **A**^t:

$$\mathbf{A}^{t} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}\mathbf{Q}\Lambda\mathbf{Q}^{-1}\cdots\mathbf{Q}^{-1}\Lambda\mathbf{Q}^{-1} = \mathbf{Q}\Lambda^{t}\mathbf{Q}^{-1},$$
$$\Lambda^{t} = \operatorname{diag}(\lambda_{1}^{t},...,\lambda_{D}^{t})$$

• If λ is an eigenvalue of **A**, then $\lambda + \alpha$ is an eigenvalue of **A** + α **I**:

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$:

• If λ is an eigenvalue of **A**, then λ^t is an eigenvalue of **A**^{*t*}:

$$\mathbf{A}^{t} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}\mathbf{Q}\Lambda\mathbf{Q}^{-1}\cdots\mathbf{Q}^{-1}\Lambda\mathbf{Q}^{-1} = \mathbf{Q}\Lambda^{t}\mathbf{Q}^{-1},$$
$$\Lambda^{t} = \operatorname{diag}(\lambda_{1}^{t},...,\lambda_{D}^{t})$$

• If λ is an eigenvalue of **A**, then $\lambda + \alpha$ is an eigenvalue of **A** + α **I**:

$$\mathbf{A} + \alpha \mathbf{I} = \mathbf{Q} \Lambda \mathbf{Q}^{-1} + \mathbf{Q} \alpha \mathbf{I} \mathbf{Q}^{-1} = \mathbf{Q} (\Lambda + \alpha \mathbf{I}) \mathbf{Q}^{-1},$$
$$\Lambda + \alpha \mathbf{I} = \text{diag}(\lambda_1 + \alpha, ..., \lambda_D + \alpha)$$

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$:

• If λ is an eigenvalue of **A**, then λ^t is an eigenvalue of **A**^t:

$$\mathbf{A}^{t} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}\mathbf{Q}\Lambda\mathbf{Q}^{-1}\cdots\mathbf{Q}^{-1}\Lambda\mathbf{Q}^{-1} = \mathbf{Q}\Lambda^{t}\mathbf{Q}^{-1},$$
$$\Lambda^{t} = \operatorname{diag}(\lambda_{1}^{t},...,\lambda_{D}^{t})$$

• If λ is an eigenvalue of **A**, then $\lambda + \alpha$ is an eigenvalue of **A** + α **I**:

$$\mathbf{A} + \alpha \mathbf{I} = \mathbf{Q} \Lambda \mathbf{Q}^{-1} + \mathbf{Q} \alpha \mathbf{I} \mathbf{Q}^{-1} = \mathbf{Q} (\Lambda + \alpha \mathbf{I}) \mathbf{Q}^{-1},$$
$$\Lambda + \alpha \mathbf{I} = \operatorname{diag}(\lambda_1 + \alpha, ..., \lambda_D + \alpha)$$

 Combine: If λ is an eigenvalue of A, then (λ + α)^t is an eigenvalue of (A + αI)^t

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$: Assume \mathbf{A} is symmetric

Consider the following Rayleigh quotient

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^{\top} \mathbf{A} \mathbf{x}}{||\mathbf{x}||_{2}^{2}}, \quad ||\mathbf{x}||_{2}^{2} = \mathbf{x}^{\top} \mathbf{x}$$

• We can show that

$$\lambda_{min}(\mathbf{A}) \leq R(\mathbf{A}, \mathbf{x}) \leq \lambda_{max}(\mathbf{A})$$
$$\Rightarrow \quad \lambda_{min}(\mathbf{A}) ||\mathbf{x}||_2^2 \leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \lambda_{max}(\mathbf{A}) ||\mathbf{x}||_2^2$$

(break)

Indications of $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1}$: Assume \mathbf{A} is symmetric

Consider the following Rayleigh quotient

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^{\top} \mathbf{A} \mathbf{x}}{||\mathbf{x}||_{2}^{2}}, \quad ||\mathbf{x}||_{2}^{2} = \mathbf{x}^{\top} \mathbf{x}$$

• We can show: $\lambda_{min}(\mathbf{A}) \leq R(\mathbf{A}, \mathbf{x}) \leq \lambda_{max}(\mathbf{A})$

$$\Rightarrow \quad \lambda_{min}(\mathbf{A})||\mathbf{x}||_2^2 \leqslant \mathbf{x}^\top \mathbf{A} \mathbf{x} \leqslant \lambda_{max}(\mathbf{A})||\mathbf{x}||_2^2$$

Gradient descent with constant step-size to find θ^* :

$$\boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*, \quad \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

• The ℓ_2 distance between θ_t and θ^* :

$$\begin{aligned} ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 &= ||(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)||_2^2 \\ &= |(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)| \end{aligned}$$

Gradient descent with constant step-size to find θ^* :

$$\boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*, \quad \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

• The ℓ_2 distance between θ_t and θ^* :

$$||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 = ||(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)||_2^2$$
$$= |(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)|$$

• Bounded distance by setting $x = \theta_0 - \theta^*$, $\mathbf{A} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t}$:

$$||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \ge \lambda_{min}((\mathbf{I} - \frac{\gamma}{\sigma^{2}}\mathbf{X}^{\top}\mathbf{X})^{2t})||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \le \lambda_{max}((\mathbf{I} - \frac{\gamma}{\sigma^{2}}\mathbf{X}^{\top}\mathbf{X})^{2t})||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$

Gradient descent with constant step-size to find θ^* :

$$\boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*, \quad \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

• The ℓ_2 distance between θ_t and θ^* :

$$||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 = ||(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)||_2^2$$
$$= |(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)|$$

• Bounded distance by setting $x = \theta_0 - \theta^*$, $\mathbf{A} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t}$:

$$||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 \ge \lambda_{min}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2)^t ||\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*||_2^2$$
$$||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 \le \lambda_{max}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2)^t ||\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*||_2^2$$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq ||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq \lambda_{max}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2})$$

,

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq ||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq \lambda_{max}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2})$$

Convergence properties in difference cases:

- 1. $\lambda_{max} < 1$: always converge
- 2. $\lambda_{min} \ge 1$: always diverge
- 3. $\lambda_{min} < 1$ but $\lambda_{max} \ge 1$: convergence depending on θ_0

Gradient descent with constant step-size to find θ^* :

$$\begin{split} \lambda_{min}^{t} || \boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*} ||_{2}^{2} &\leq || \boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*} ||_{2}^{2} \leq \lambda_{max}^{t} || \boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*} ||_{2}^{2} \\ \lambda_{min} &:= \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \\ \text{Deriving the eigenvalues } \lambda_{min}, \lambda_{max} : \end{split}$$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^t || \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* ||_2^2 \leqslant || \boldsymbol{\theta}_t - \boldsymbol{\theta}^* ||_2^2 \leqslant \lambda_{max}^t || \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* ||_2^2$$

$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\mathsf{T}} \mathbf{X})^2) \ge 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\mathsf{T}} \mathbf{X})^2)$$

Deriving the eigenvalues λ_{min} , λ_{max} :

• If λ is an eigenvalue of $\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}$, then λ^2 is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^t || \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* ||_2^2 \leqslant || \boldsymbol{\theta}_t - \boldsymbol{\theta}^* ||_2^2 \leqslant \lambda_{max}^t || \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* ||_2^2$$

$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\mathsf{T}} \mathbf{X})^2) \ge 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\mathsf{T}} \mathbf{X})^2)$$

Deriving the eigenvalues λ_{min} , λ_{max} :

- If λ is an eigenvalue of $\mathbf{I} \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}$, then λ^2 is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2$
- If λ is an eigenvalue of $\mathbf{X}^{\top}\mathbf{X}$, then $1 - \frac{\gamma\lambda}{\sigma^2}$ is an eigenvalue of $\mathbf{I} - \frac{\gamma}{\sigma^2}\mathbf{X}^{\top}\mathbf{X}$:

$$\mathbf{X}^{\mathsf{T}} \mathbf{X} \boldsymbol{q} = \lambda \boldsymbol{q} \quad \Leftrightarrow \quad (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\mathsf{T}} \mathbf{X}) \boldsymbol{q} = (1 - \frac{\gamma \lambda}{\sigma^2}) \boldsymbol{q}$$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq ||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq \lambda_{max}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2})$$

• If λ is an eigenvalue of $\mathbf{X}^{\top}\mathbf{X}$, then $(1 - \frac{\gamma\lambda}{\sigma^2})^2$ is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2}\mathbf{X}^{\top}\mathbf{X})^2$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq ||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq \lambda_{max}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2})$$

- If λ is an eigenvalue of $\mathbf{X}^{\top}\mathbf{X}$, then $(1 - \frac{\gamma\lambda}{\sigma^2})^2$ is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2}\mathbf{X}^{\top}\mathbf{X})^2$
- $\mathbf{X}^{\top}\mathbf{X}$ is positive semi-definite $\Rightarrow \lambda \ge 0$

Gradient descent with constant step-size to find θ^* :

$$\lambda_{min}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq ||\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{*}||_{2}^{2} \leq \lambda_{max}^{t} ||\boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{*}||_{2}^{2}$$
$$\lambda_{min} := \lambda_{min} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2}) \geq 0, \quad \lambda_{max} := \lambda_{max} ((\mathbf{I} - \frac{\gamma}{\sigma^{2}} \mathbf{X}^{\top} \mathbf{X})^{2})$$

- If λ is an eigenvalue of $\mathbf{X}^{\top}\mathbf{X}$, then $(1 - \frac{\gamma\lambda}{\sigma^2})^2$ is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2}\mathbf{X}^{\top}\mathbf{X})^2$
- $\mathbf{X}^{\top}\mathbf{X}$ is positive semi-definite $\Rightarrow \lambda \ge 0$
- Ensuring convergence: we want $\lambda_{max} = \max(1 \frac{\gamma\lambda}{\sigma^2})^2 < 1$

$$\Rightarrow \gamma < \frac{2\sigma^2}{\lambda_{max}(\mathbf{X}^{\top}\mathbf{X})}$$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes?

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky.

1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top} \mathbf{X})$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky.

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top} \mathbf{X})$
- 2. You choose an initialisation θ_0
 - Consider eigen decomposition of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ as $\{(\lambda_d, \mathbf{q}_d)\}$ with

 $\lambda_1 \ge ... \ge \lambda_D \ge 0$ (therefore $\lambda_{max}(\mathbf{X}^{\top}\mathbf{X}) = \lambda_1$)

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky.

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top} \mathbf{X})$
- 2. You choose an initialisation θ_0
 - Consider eigen decomposition of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ as $\{(\lambda_d, q_d)\}$ with

 $\lambda_1 \ge ... \ge \lambda_D \ge 0$ (therefore $\lambda_{max}(\mathbf{X}^{\top}\mathbf{X}) = \lambda_1$)

3. Lucky case:

• now assume it happens to be that $\theta_0 = \theta^* + \alpha q_d, d > 1$

$$\begin{aligned} ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 &= \alpha^2 |\boldsymbol{q}_d^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} \boldsymbol{q}_d| = \alpha^2 (1 - \frac{\gamma \lambda_d}{\sigma^2})^{2t} \\ \bullet \ \gamma &< 2\sigma^2 / \lambda_d \quad \Rightarrow \quad ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 \to 0 \end{aligned}$$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky.

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top} \mathbf{X})$
- 2. You choose an initialisation θ_0
 - Consider eigen decomposition of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ as $\{(\lambda_d, q_d)\}$ with

 $\lambda_1 \ge ... \ge \lambda_D \ge 0$ (therefore $\lambda_{max}(\mathbf{X}^{\top}\mathbf{X}) = \lambda_1$)

3. Lucky case:

• now assume it happens to be that $\theta_0 = \theta^* + \alpha q_d, d > 1$

$$\begin{aligned} ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 &= \alpha^2 |\boldsymbol{q}_d^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} \boldsymbol{q}_d| = \alpha^2 (1 - \frac{\gamma \lambda_d}{\sigma^2})^{2t} \\ &\simeq \gamma < 2\sigma^2 / \lambda_d \quad \Rightarrow \quad ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 \to 0 \end{aligned}$$

4. **Unlucky case**: $\gamma \ge 2\sigma^2/\lambda_d \Rightarrow \text{divergence}$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky.

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top} \mathbf{X})$
- 2. You choose an initialisation θ_0
 - Consider eigen decomposition of $\mathbf{X}^{\top}\mathbf{X}$ as $\{(\lambda_d, \mathbf{q}_d)\}$ with

 $\lambda_1 \ge ... \ge \lambda_D \ge 0$ (therefore $\lambda_{max}(\mathbf{X}^{\top}\mathbf{X}) = \lambda_1$)

3. Lucky case:

• now assume it happens to be that $\theta_0 = \theta^* + \alpha q_d, d > 1$

$$||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 = \alpha^2 |\boldsymbol{q}_d^\top (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t} \boldsymbol{q}_d| = \alpha^2 (1 - \frac{\gamma \lambda_d}{\sigma^2})^{2t}$$

•
$$\gamma < 2\sigma^2/\lambda_d \quad \Rightarrow \quad ||\boldsymbol{\theta}_t - \boldsymbol{\theta}^*||_2^2 \to 0$$

4. Unlucky case: $\gamma \ge 2\sigma^2/\lambda_d \implies \text{divergence}$

Constant step-size GD diverges if $\gamma \ge 2\sigma^2 / \lambda_{min}(\mathbf{X}^{\top}\mathbf{X})$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2/\lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$
- 2. You randomly initialise the parameter θ_0

• write
$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \boldsymbol{q}_d$$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$

2. You randomly initialise the parameter θ_0

• write
$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \boldsymbol{q}_d$$

3. Evolution of the iterative update: $\boldsymbol{\theta}_t - \boldsymbol{\theta}^* = \sum_{d=1}^{D} \alpha_d (1 - \frac{\gamma \lambda_d}{\sigma^2})^t \boldsymbol{q}_d$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

- 1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$
- 2. You randomly initialise the parameter θ_0

• write
$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \boldsymbol{q}_d$$

3. Evolution of the iterative update: $\theta_t - \theta^* = \sum_{d=1}^{D} \alpha_d (1 - \frac{\gamma \lambda_d}{\sigma^2})^t q_d$

• For direction
$$q_d$$
 with $\gamma < 2\sigma^2/\lambda_d$: $(1 - \frac{\gamma\lambda_d}{\sigma^2})^t q_d \to 0$

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$

2. You randomly initialise the parameter θ_0

• write
$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \boldsymbol{q}_d$$

3. Evolution of the iterative update: $\theta_t - \theta^* = \sum_{d=1}^{D} \alpha_d (1 - \frac{\gamma \lambda_d}{\sigma^2})^t q_d$

- For direction q_d with $\gamma < 2\sigma^2/\lambda_d$: $(1 \frac{\gamma\lambda_d}{\sigma^2})^t q_d \rightarrow \mathbf{0}$
- For other directions $\gamma \ge 2\sigma^2/\lambda_d$: $(1 \frac{\gamma\lambda_d}{\sigma^2})^t q_d$ diverges
- GD diverges unless $\alpha_d = 0$ for those *d* with $\gamma \ge 2\sigma^2/\lambda_d$.

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ **Q**: Can we use larger step-sizes? **A**: Yes if you are lucky, and no if you are unlucky. **Caveat:** you are unlikely to be that lucky...

1. You choose a step-size $\gamma > 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$

2. You randomly initialise the parameter θ_0

• write
$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \boldsymbol{q}_d$$

3. Evolution of the iterative update: $\theta_t - \theta^* = \sum_{d=1}^{D} \alpha_d (1 - \frac{\gamma \lambda_d}{\sigma^2})^t q_d$

- For direction q_d with $\gamma < 2\sigma^2/\lambda_d$: $(1 \frac{\gamma\lambda_d}{\sigma^2})^t q_d \rightarrow \mathbf{0}$
- For other directions $\gamma \ge 2\sigma^2/\lambda_d$: $(1 \frac{\gamma\lambda_d}{\sigma^2})^t q_d$ diverges
- GD diverges unless $\alpha_d = 0$ for those *d* with $\gamma \ge 2\sigma^2/\lambda_d$.
- 4. For d > 1, span{ q_d , ..., q_D } has measure 0 in \mathbb{R}^D
 - unlikely to make $\alpha_d = 0$ for some *d* with random initialisation

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ If you want to test your luck: choose $\gamma \in \left[\frac{2\sigma^2}{\lambda_{max}(\mathbf{X}^{\top}\mathbf{X})}, \frac{2\sigma^2}{\lambda_{min}(\mathbf{X}^{\top}\mathbf{X})}\right)$ Is my choice of γ robust to initialisation of $\boldsymbol{\theta}_0$?

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ If you want to test your luck: choose $\gamma \in \left[\frac{2\sigma^2}{\lambda_{max}(\mathbf{X}^{\top}\mathbf{X})}, \frac{2\sigma^2}{\lambda_{min}(\mathbf{X}^{\top}\mathbf{X})}\right)$ Is my choice of γ robust to initialisation of θ_0 ?

Depending on the condition number:



Need careful choice of step-sizes if the loss is "very stretched"

To ensure convergence at **any** initialisation: $\gamma < 2\sigma^2 / \lambda_{max}(\mathbf{X}^{\top}\mathbf{X})$ If you want to test your luck: choose $\gamma \in [\frac{2\sigma^2}{\lambda_{max}(\mathbf{X}^{\top}\mathbf{X})}, \frac{2\sigma^2}{\lambda_{min}(\mathbf{X}^{\top}\mathbf{X})})$ Is my choice of γ robust to initialisation of $\boldsymbol{\theta}_0$?

Depending on the condition number:



Need careful choice of step-sizes if the loss is "very stretched"

• Note:
$$\kappa(\mathbf{X}^{\top}\mathbf{X}) = \kappa(\mathbf{X})^2 = \frac{\sigma_{max}(\mathbf{X})}{\sigma_{min}(\mathbf{X})}$$

Choosing step-size: general case

In general the loss function is non-quadratic nor convex:



https://distill.pub/2017/momentum/

Gradient Descent

Yingzhen Li @Imperial College London, October 26, 2021
In general the loss function is non-quadratic nor convex:



Local quadratic approximation when $\theta_t \approx \theta^*$:

- locally approximate $L(\boldsymbol{\theta}_t) \approx L(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta}_t \boldsymbol{\theta}^*)^\top \nabla^2 L(\boldsymbol{\theta}^*)(\boldsymbol{\theta}_t \boldsymbol{\theta}^*)$ (in linear regression $\nabla^2 L(\boldsymbol{\theta}) \propto \mathbf{X}^\top \mathbf{X}$)
- $\kappa(\nabla^2 L)$ can tell whether the loss is "locally stretched"

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Let's see what happens for different step-sizes.



Image shows:

- Path of θ_t from Gradient Descent
- Constant step size $\gamma_t = \gamma$

Choosing step-size: summary

Summary on choosing step size:

- too small: slow convergence
- too large: divergence
- just right: depends on problem (often: trial and error)

Choosing step-size: summary

Summary on choosing step size:

- too small: slow convergence
- too large: divergence
- just right: depends on problem (often: trial and error)

Rule of thumb: Start from a relatively large step size, decrease step size as getting closer to a (local) optimum.