# Optimization theory

# Probability theory

univariate

1. Random variable: sample space $\Omega$., $P$ ... probability measure

$\bigstar$ $p$ ... density function

$\Phi$ ... cdf.

$\triangleright$ Expected value:
- $E[X]$, $\mu_X$
- LLN.
- linear : $E\left[\begin{smallmatrix} \text{linear} \\ \text{combination} \end{smallmatrix}\right] = \begin{smallmatrix} \text{linear} \\ \text{comb.} \end{smallmatrix}\left(E[\dots]\right)$.

$\triangleright$ Variance:
- $\text{Var}[X]$, $\sigma_X^2$ ... $\sigma_X$ ... standard deviation.
- Not linear $\text{Var}\left[\begin{smallmatrix} \text{linear} \\ \text{comb} \end{smallmatrix}\right] = \dots \alpha^2$, $\text{Cov}$

2. Multivariate R.V. (vectors)
  ⊢ Marginal distribution
  ⊢ Conditional RV ⟶ Bayes' theorem.
  ⊢ Independence , correlation.

3. Parametric families
  ⊢ Bernoulli              ⊢ Uniform
  ⊢ Binomial.              ⊢ Gaussian (CLT).
  ⊢ Poison
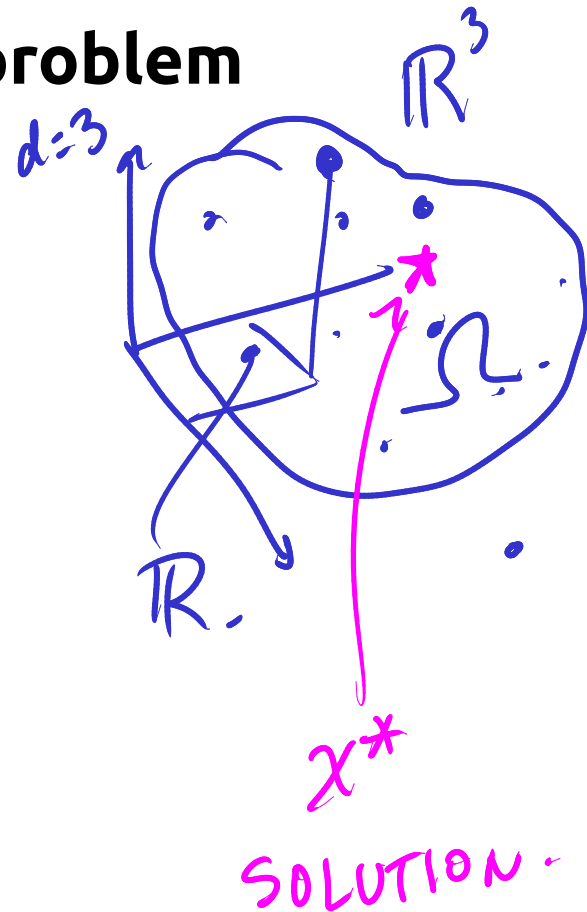  ⊢ Exponential

# Components of an optimization problem

1. The **decision vector** : $x \in \mathbb{R}^d$

   dimension

2. The **search set** or **feasible set** : $\Omega \subseteq \mathbb{R}^d$

   not the sample space !

3. The **objective function** : $J : \Omega \to \mathbb{R}$

$\mathbb{R}^3$

$d=3$

$\Omega$

$\mathbb{R}$.

$x^*$

SOLUTION.

**Notation.**

minimize$_x$ $J(x)$

subject to: $x \in \Omega$

MAXIMIZE$_x$ $J(x)$

s.t. $x \in \Omega$

$\Longleftrightarrow$ minimize$_x$ $-J(x)$

s.t. $x \in \Omega$

**Constraint specification**

minimize$_x$ $J(x)$

subject to:
$f_i(x) = 0 \quad i = 1 \dots n$ ... equality
$g_j(x) \leq 0 \quad j = 1 \dots m$ ... inequality.

$x^* = \underset{x}{\text{argmin}} \quad J(x)$

subject to:
$f_i(x) = 0 \quad i = 1 \dots n$
$g_j(x) \leq 0 \quad j = 1 \dots m$

$d = 1$

$J(x)$

$\Omega$

4 — $\Omega$ — 6 — $\mathbb{R}$

$x \geq 4$ $x \leq 6$

$x_2 \leq g(x_1)$

$x_2 - g(x_1) \leq 0$

$f(x_1, x_2) = 0$

$g(x_1)$

$x_2$

$x_1$

$\Omega$

$d = 1$

$J(x)$

● ... stationary points : $J'(x) = 0$

● ... boundary points.

$J(x^*)$

$\Omega$  $x^*$

$x$

Solution $x^*$ is a stationary point.

$d = 2$

$J(x)$.

$x_2$

$x_1$

$\Omega$ ... $x^*$

local minimum.

stationary

local maximum

$\Omega$

local minimum.

$x^*$

solution.
local minimum.

**Example:** Design a rabbit enclosure



Length $\ell$ of fence material

$d = 2$

$x = (\delta, w)$.

$J(\delta, w) = -\delta w$.

$\mathbb{R}^2$

Lines of constant $w\delta$

$\Omega$

$\ell/2$

$w^* = \ell/4$

$\delta^* = \ell/4$    $\ell/2$    $\delta$

minimize    $-\delta w$.
$\delta, w$

s.t.    $\delta \geq 0$
        $w \geq 0$   ] 2 inequality constraints
        $2\delta + 2w = \ell$ ] 1 equality constraint.

# Global vs. local solutions

small ball
centered on $x^+$

- Global solution:  $J(x^*) \leq J(x) \qquad \forall\, x \in \Omega$

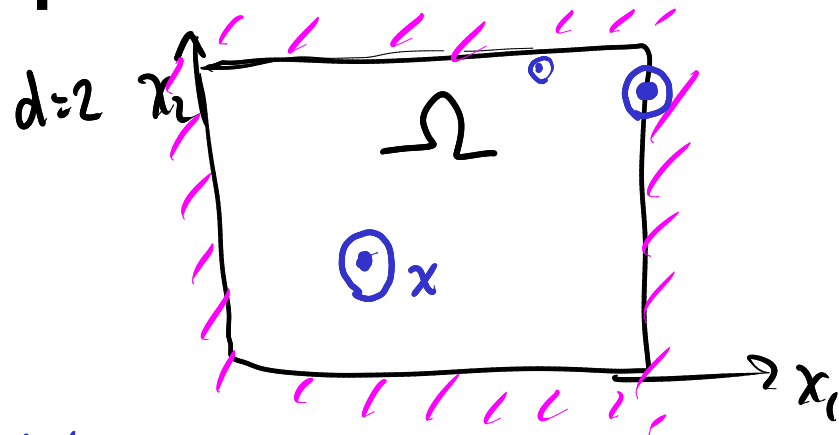- Local solution :  $J(x^+) \leq J(x) \qquad \forall\, x \in \Omega \cap B_\epsilon(x^+)$



$d = 1$

$J(x)$

$\Omega$

$x^+$

# Types of feasible points

- Interior vs. non-interior points $\dots \Omega$

There is a ball that is contained in $\Omega$

a.k.a. **boundary**
There is no ball

$d=2$  $x_2$

$\Omega$

$\odot x$

$\rightarrow x_1$

- Differentiable vs. non-differentiable points

everything else.

discontinuos (e.g $x_d$)
non-differentiable (e.g. $x_c$).

$d=1$

$J(x)$

$x_d$  $x_c$

- Stationary vs. non-stationary points

$\nabla J = 0$

$\nabla J \neq 0$

(a)  (b)  (c)  (d)  (e)

$J' = 0$ inflection point

boundary local min.

Stationary points : local minima : $J''$ (Hessian) $> 0$.
local maxima $J'' < 0$.
inflection points. otherwise : $J'' = 0$.

# First order optimality condition

$x$ is a differentiable, interior, local solution $\Rightarrow$ $x$ is stationary

All points $x$:

→1. Differentiable and interior .... All solutions in 1st category are stationary.

→2. Non-differentiable

3. Non-interior

All local solutions are stationary, Non-diff., or Non-interior

boundary.

## Corollary.

If $J$ is smooth. and no constraints.

$\Longrightarrow$ all local solutions are amongst the stationary points.

$$\underset{x,y}{\text{minimize}} \quad 100(y-x^2)^2 + (1-x)^2$$



1) $\dfrac{\partial J}{\partial x} = 200(y-x^2)(-2x) + 2(1-x)(-1) = 0$

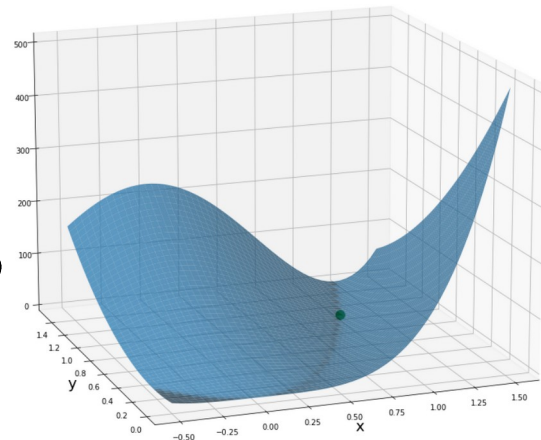$\dfrac{\partial J}{\partial y} = 200(y-x^2) = 0. \longrightarrow \quad y = x^2$

$\hookrightarrow x = 1 \implies y = 1$

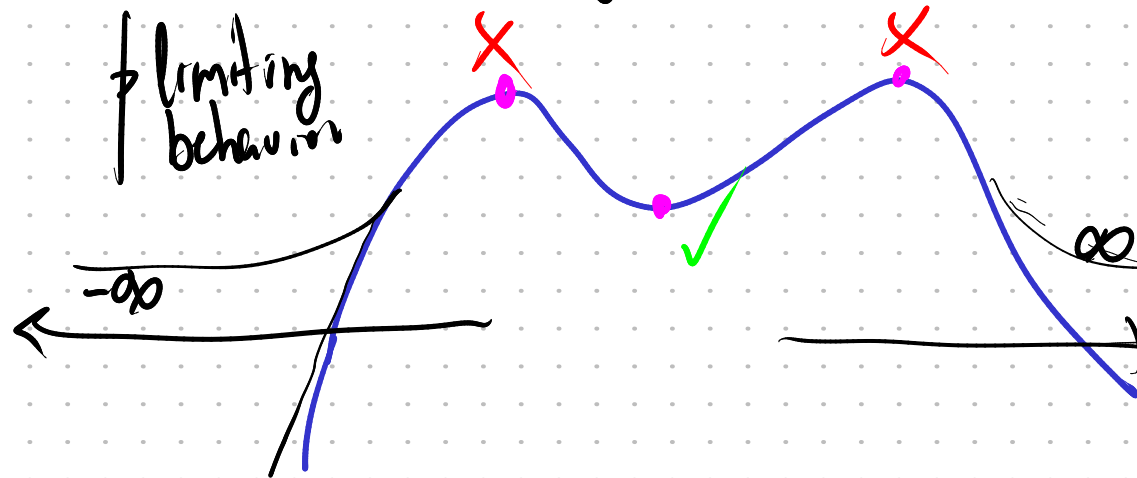**Unique stationary point.**

$(x,y) = (1,1)$

2) Verify local minimum....

↳ Check the 2nd derivative.

↳ Plot it.

---

local minimum $\neq\!\!\Rightarrow$ global minimum.
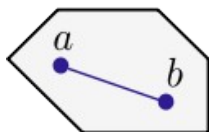
3. Check that a global solution exists.

↳ limiting behavior

No solution here.

$-\infty$

$\infty$

$\longrightarrow J(x^*) \leq J(x) \ \forall x \in \Omega.$

# Convex optimization problems

convex function

$$\underset{x}{\text{minimize}} \quad J(x)$$

subject to: $\quad x \in \Omega$

convex set

epigraph $J$ is covex set

$J$

$x$
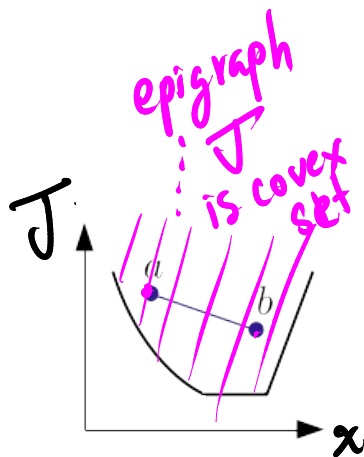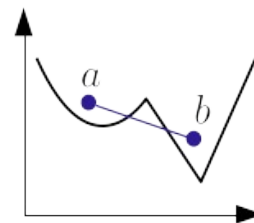
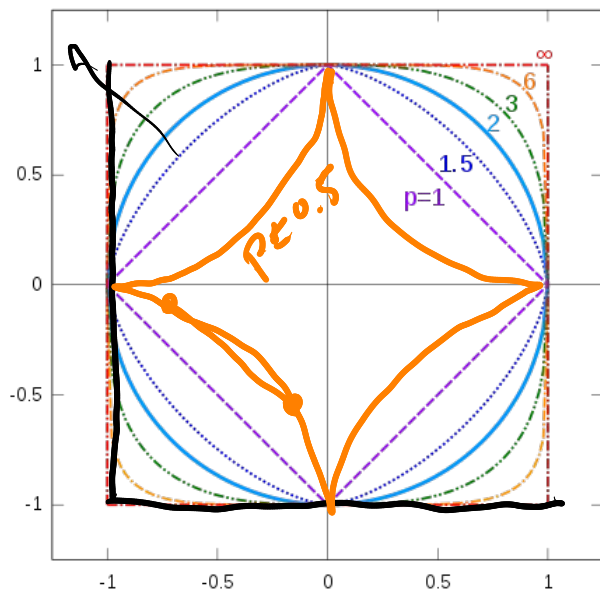Convex set          Non-convex set          Convex function          Non-convex function

# Properties of convex problems

1) For convex problems, every local solution is a global solution.

2) For unconstrained convex problems with continuously differentiable cost function, every stationary point is a global solution.

# Examples of convex sets

- p-norm ball : $\{x \in \mathbb{R}^{\textcircled{d}}: \|x\|_p \leq \textcircled{r}\}$  with  $\|x\|_p = \left(\sum_{i=1}^{d} |x_i|^p\right)^{1/p}$

$$P \geq 1$$



2 norm ball :

$$\chi_1^2 + \chi_2^2 + \chi_3^2 \leq 100 \qquad / r^2$$

3 norm  $\sqrt[3]{|\chi_1|^3 + |\chi_d|^3} \leq \#$

$\infty$ norm  $\max\left(|\chi_1|, |\chi_2|, \ldots, |\chi_d|\right) \leq \#$

# Examples of convex sets

- Affine equality constraints: $Ax = b$

Hyperplane in $\mathbb{R}^d$

convex function.

- Convex inequality constraints: $g(x) \leq 0$

$d=2$

$g(x)$

convex

convex.

$x_2$

$x_1$

# Examples of convex functions

- Affine functions: $J(x) = a^T x + b$ $\quad a \in \mathbb{R}^d, b \in \mathbb{R}$

- p-norms: $J(x) = \|x\|_p$ $\quad p \geq 1$

- Function composition: $J(x) = g(h(x))$ $\quad$ $h$ convex

  $g$ ~~non-decreasing~~

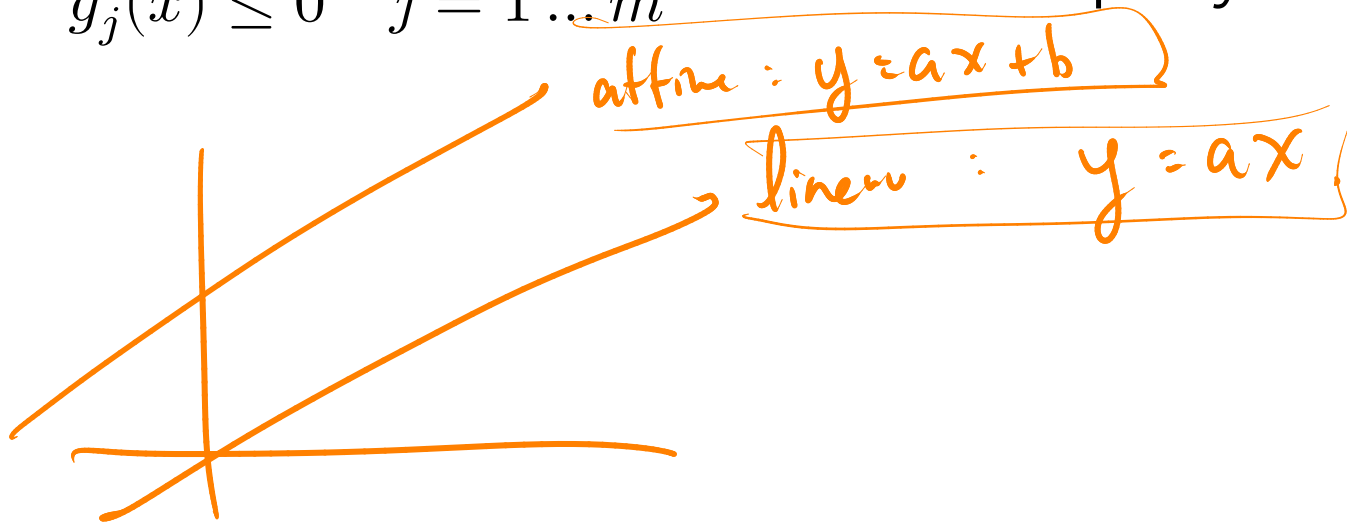  affine.

# Convex optimization problems

minimize $J(x)$ ... convex cost function
$x$

subject to: $f_i(x) = 0 \quad i = 1 \ldots n$ ... affine equality constraints

$g_j(x) \leq 0 \quad j = 1 \ldots m$ ... convex inequality constraints

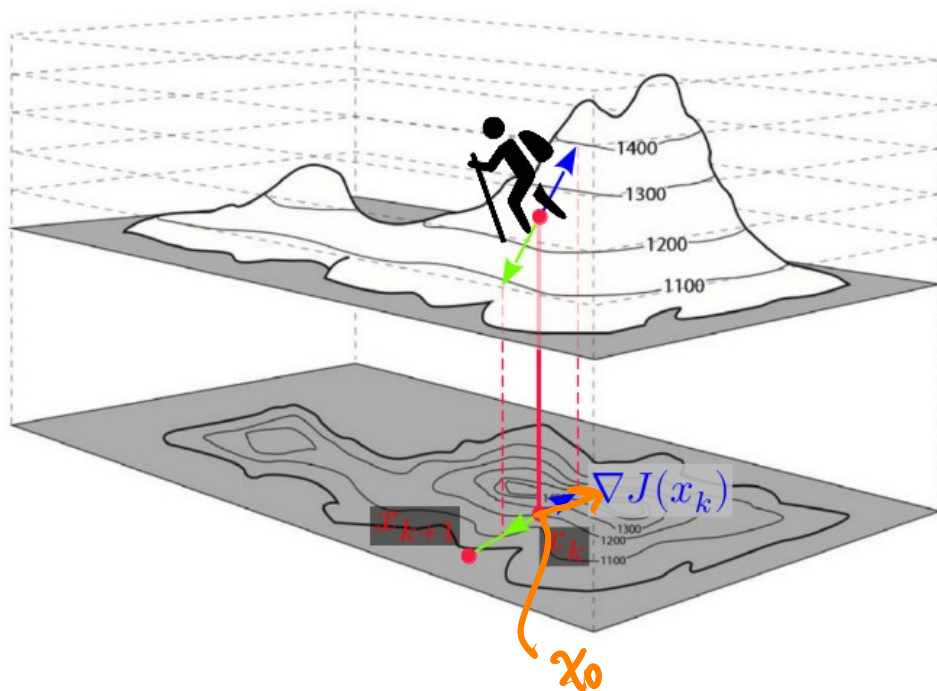affine: $y = ax + b$

linear: $y = ax$

# Gradient descent ... finds local minima.

$$x^* = \operatorname*{argmin}_x J(x)$$

k... step counter.



0. Initialize: $x_0 \in \Omega$, $k = 0$

1. Loop until convergence:
   - $x_{k+1} = x_k - \gamma \nabla J(x_k)$
   - $k \leftarrow k + 1$
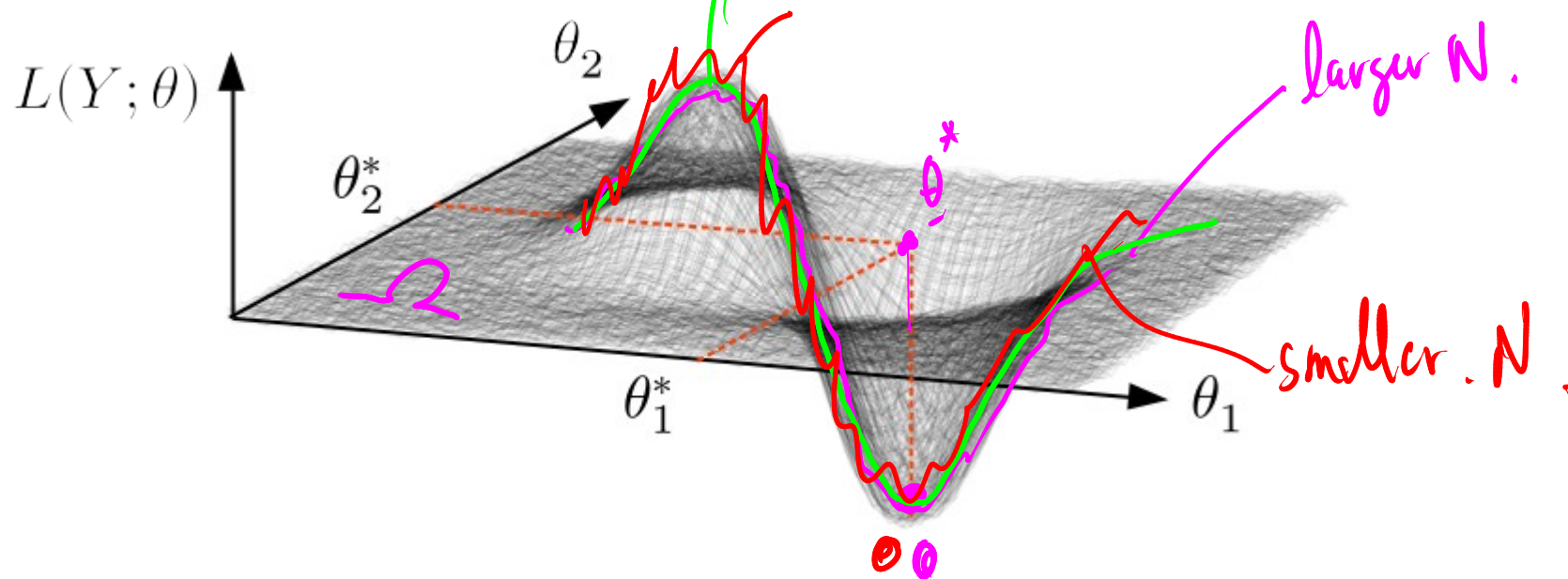
stepsize

$x_0$

$\nabla J(x_k)$

# Stochastic optimization

$$\underline{\theta} = (\theta_1, \ldots, \theta_D).$$

$\theta \ldots$ scalar.

Example:  $\underline{\theta} = (\theta_1, \theta_2)$

$$\theta^* = \operatorname*{argmin}_{\underline{\theta}} E[L(Y \, ; \, \underline{\theta})]$$

loss.

N is small

larger N.

smaller N.

Full problem : $\quad \underline{\theta}^* = \underset{\underline{\theta}}{\mathsf{argmin}}\ E[L(Y\ ;\ \underline{\theta})]$

$\mathcal{D} = \{y_i\}_N \overset{iid}{\sim} Y$

Approximate expectation : $\quad E[L(Y\ ;\ \underline{\theta})] \approx \dfrac{1}{N}\sum_{i=1}^{N} L(y_i\ ;\ \underline{\theta})$

$J(\underline{\theta})$

Approximate problem : $\quad \underline{\theta}^* = \underset{\underline{\theta}}{\mathsf{argmin}}\ \sum_{i=1}^{N} L(\underline{\theta}\ ;\ y_i)$

Gradient descent : $\quad \underline{\theta}_{k+1} = \underline{\theta}_k - \gamma \nabla_{\underline{\theta}}\left(\sum_{i=1}^{N} L(\underline{\theta}_k\ ;\ y_i)\right)$

$\qquad\qquad = \underline{\theta}_k - \gamma \sum_{i=1}^{N} \nabla_{\underline{\theta}} L(\underline{\theta}_k\ ;\ y_i)$
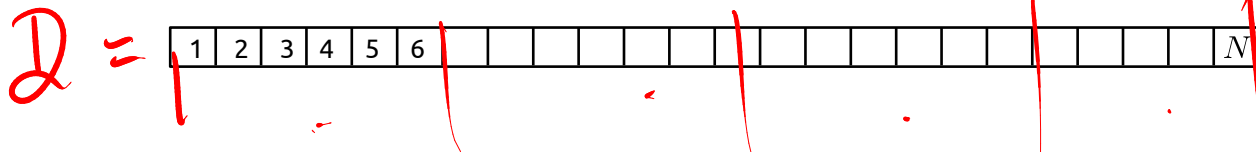
$\mathcal{D}.$

# Stochastic gradient descent (SGD)

$$\underline{\theta}_{k+1} = \underline{\theta}_k - \gamma \sum_{i \in \mathcal{B}} \nabla_{\underline{\theta}} L(\underline{\theta}\,;\, y_i)$$

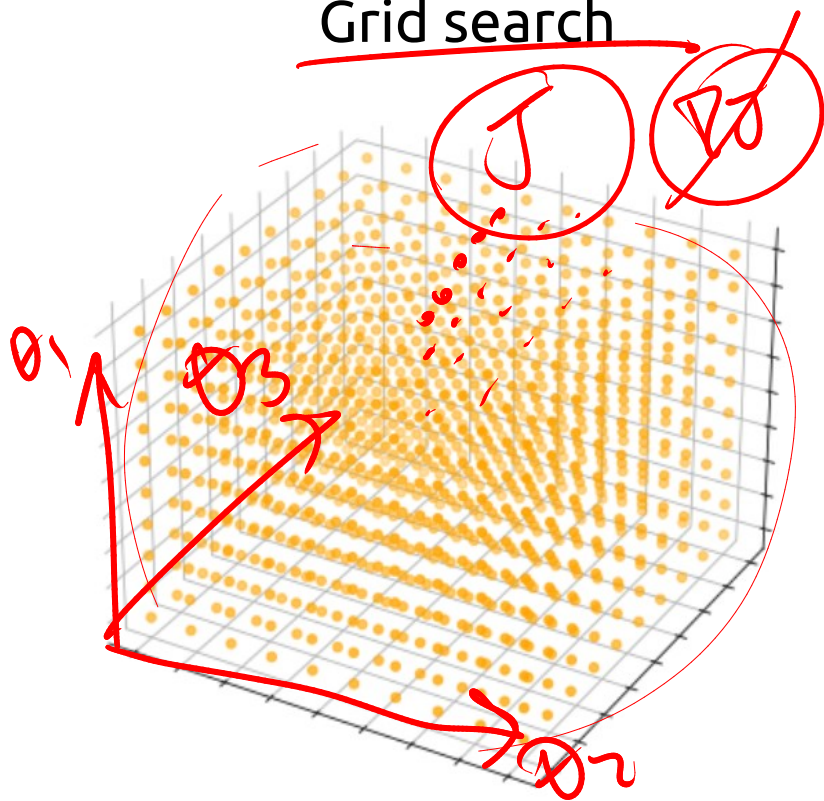$\mathcal{B} \subseteq \{1 \dots N\}$

a *batch* of samples

$\mathcal{D} =$

| 1 | 2 | 3 | 4 | 5 | 6 | | | | | | | | | | | | | | | | | | | N |

epoch.

- $|\mathcal{B}| = 1$ (pure SGD)
  - $\mathcal{B}$ chosen *with replacement*
  - $\mathcal{B}$ chosen *without replacement*
- $1 < |\mathcal{B}| < N$ (minibatch)
  - partition
  - shuffle / split
- $|\mathcal{B}| = N$ (regular gradient descent)

# Gradient-less optimization



Grid search

Genetic algorithms