**Applied Machine Learning (AML)**

Optimisation

Oisin Mac Aodha • Siddharth N.

---

**Optimisation**

---

# Core Questions

- What task am I trying to solve?
- How should I model the problem?
- How should I represent my data?
- How can I estimate the parameters of my model?
- How should I measure the performance of my model?

---

# Why Optimisation?

**Main idea:** "learning" $\rightarrow$ continuous optimisation

- Linear regression
- Logistic regression
- Neural networks

...

**Maximum Likelihood**

$$\ell(\boldsymbol{w}) = \log p(\boldsymbol{x}_1, y_1, \boldsymbol{x}_2, y_2, \ldots, \boldsymbol{x}_N, y_N | \boldsymbol{w})$$

$$= \log \prod_{i=1}^{N} p(\boldsymbol{x}_i, y_i, | \boldsymbol{w})$$

$$= \sum_{i=1}^{N} \log p(\boldsymbol{x}_i, y_i, | \boldsymbol{w})$$

E.g. $\text{NLL}(\boldsymbol{w}) = -\ell(\boldsymbol{w})$
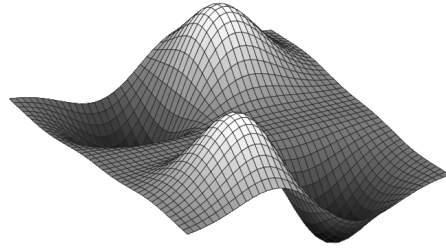
## Why Optimisation?

**Result:** An "error function" $\mathcal{L}(w)$ to minimise

### Error function

- For fixed data $\mathcal{D}$, every setting of weights results in some error
- Learning $\equiv$ descending error surface
- When data is iid

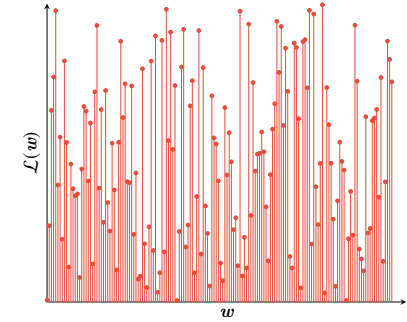$$\mathcal{L}(w) = \sum_i \mathcal{L}^i(w)$$

for each data point

---

## Role of Smoothness

### Unconstrained

- minimisation impossible
- …can only search through all $w$

### Constrained/Continuous

- $\mathcal{L}(w)$ provides information about $\mathcal{L}$ at nearby values

---

## Role of Smoothness

### Unconstrained

- minimisation impossible
- …can only search through all $w$

### Constrained/Continuous

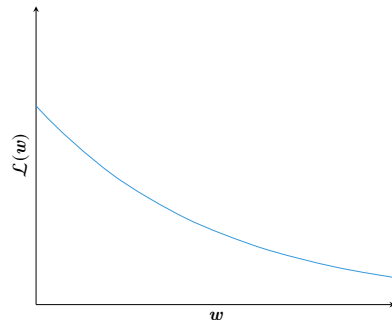- $\mathcal{L}(w)$ provides information about $\mathcal{L}$ at nearby values

---

## Role of Derivatives

**Check:** perturb $w_i$ keeping all else the same; does error $\uparrow / \downarrow$ ?

### Calculus

- for differentiable $\mathcal{L}$, partial derivatives $\frac{\partial \mathcal{L}}{\partial w_i}$
- vector of partial derivatives $\equiv$ gradient of the error

$$\nabla_w \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \ldots, \frac{\partial \mathcal{L}}{\partial w_N} \right)$$

- direction of steepest error *ascent*

  ($-\nabla_w \mathcal{L}$ for *descent*)

### Key Challenges

- Computing $\nabla_w \mathcal{L}$ efficiently
- Minimising error with gradient
- Location of minimiser

# Optimisation Algorithms

## Optimisation

---

**General Optimisation Problem**

$$\min_{w} \mathcal{L}(w)$$

### Components

- procedure to compute $\mathcal{L}(w)$
- procedure to compute $\nabla_w \mathcal{L}(w)$

### Aside

- some others don't use gradients
- some use higher-order gradients
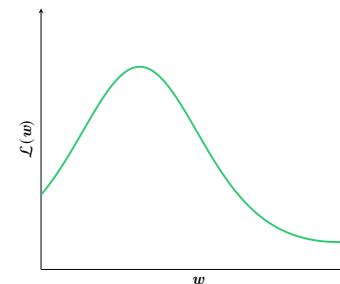
  …not covered here

---

# Basic Optimisation Algorithm

**Require:** stopping error $\epsilon$, step size $\eta$
1: $w \leftarrow$ initialisation
2: **while** $\mathcal{L}(w) > \epsilon$ **do**
3:     calculate $g = \nabla_w \mathcal{L}(w)$
4:     compute direction $d$ from $w, \mathcal{L}(w), g$
5:     $w \leftarrow w - \eta d$
6: **return** $w$

$$w_0, \quad w_1, \quad w_2, \quad \ldots$$
$$\mathcal{L}(w_0), \quad \mathcal{L}(w_1), \quad \mathcal{L}(w_2), \quad \ldots$$
$$\nabla_w \mathcal{L}(w_0), \quad \nabla_w \mathcal{L}(w_1), \quad \nabla_w \mathcal{L}(w_2), \quad \ldots$$

---

# Choosing a direction

**Simple choice:** $d = \nabla_w \mathcal{L}$    locally steepest direction

Recall (multi-variate) Taylor theorem:
for $w \in \mathbb{R}^N$ and perturbation $\delta \in \mathbb{R}^N$ such that $a = w - \delta$

$$\mathcal{L}(w) \approx \mathcal{L}(a) + \nabla_w \mathcal{L}(a)^\top \delta + \frac{1}{2} \delta^\top \nabla_w^2 \mathcal{L}(a)\delta + \ldots$$

$$\approx \mathcal{L}(a) + \nabla_w \mathcal{L}(a)^\top \delta$$

(dropping higher order terms for small $\delta$)

$$\therefore \mathcal{L}(a + \delta) = \mathcal{L}(a) + \nabla_w \mathcal{L}(a)^\top \delta$$

which is minimised at $\delta = -\eta \nabla_w \mathcal{L}(a)$ as

$$\mathcal{L}(a - \eta \nabla_w \mathcal{L}(a)) = \mathcal{L}(a) - \eta \|\nabla_w \mathcal{L}(a)\|^2$$

$$\implies \mathcal{L}(a - \eta \nabla_w \mathcal{L}(a)) \leq \mathcal{L}(a) \quad \text{(for } \eta > 0)$$

Taking a step along $\delta$ cannot increase value "locally"

# Gradient Descent

| Generic Optimisation | Gradient Descent |
|---|---|
| **Require:** stopping error $\epsilon$, step size $\eta$ | **Require:** stopping error $\epsilon$, step size $\eta$ |
| 1: $w \leftarrow$ initialisation | 1: $w \leftarrow$ initialisation |
| 2: **while** $\mathcal{L}(w) > \epsilon$ **do** | 2: **while** $\mathcal{L}(w) > \epsilon$ **do** |
| 3:      calculate $g = \nabla_w \mathcal{L}(w)$ | 3:      calculate $g = \nabla_w \mathcal{L}(w)$ |
| 4:      compute direction $d$ from $w, \mathcal{L}(w), g$ | 4:      compute direction $d = g$ |
| 5:      $w \leftarrow w - \eta d$ | 5:      $w \leftarrow w - \eta g$ |
| 6: **return** $w$ | 6: **return** $w$ |

---

# Effect of Step Size ($\eta$)

## Step Size ($\eta$)

Sometimes called *learning rate*

- $\eta > 0$
- $\eta$ too small $\rightarrow$ too slow

Example: Minimise $\mathcal{L}(w) = w^2$



Take $\eta = 0.1$

$$w_0 = 1.0$$
$$w_1 = w_0 - 0.1 \cdot 2w_0 = 0.8$$
$$w_2 = w_1 - 0.1 \cdot 2w_0 = 0.64$$
$$\vdots$$
$$w_2 = w_1 - 0.1 \cdot 2w_0 = 0.512$$
$$w_{25} = 0.0047$$

---

# Effect of Step Size ($\eta$)

## Step Size ($\eta$)

Sometimes called *learning rate*

- $\eta > 0$
- $\eta$ too small $\rightarrow$ too slow
- $\eta$ too large $\rightarrow$ instability

Example: Minimise $\mathcal{L}(w) = w^2$



Take $\eta = 1.1$

$$w_0 = 1.0$$
$$w_1 = w_0 - 1.1 \cdot 2w_0 = -1.2$$
$$w_2 = w_1 - 1.1 \cdot 2w_0 = 1.44$$
$$\vdots$$
$$w_2 = w_1 - 1.1 \cdot 2w_0 = -1.72$$
$$w_{25} = 79.50$$

## Heuristic for step size ($\eta$)

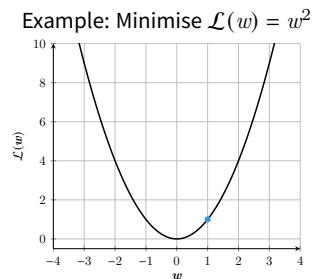**Require:** stopping error $\epsilon$, step size $\eta$

1:   $w \leftarrow$ initialisation
2:   **while** $\mathcal{L}(w) > \epsilon$ **do**
3:       compute $g = \nabla_w \mathcal{L}(w)$
4:       compute direction $d$ from $w, \mathcal{L}(w), g$
5:       $\ell_- = \mathcal{L}(w)$                 ▷ error before update
6:       $w \leftarrow w - \eta\, d$
7:       $\ell_+ = \mathcal{L}(w)$                 ▷ error after update
8:       **if** $\ell_+ \geq \ell_-$ **then**         ▷ if error increases
9:          $\eta \leftarrow \eta/2$; **revert** $w$     ▷ reduce step size
10:      **else**                        ▷ if error decreases
11:         $\eta \leftarrow 1.1\eta$           ▷ speed up *slightly*
12: **return** $w$

**Gradient Descent**

----

**Stochastic Gradient Descent**

---

## Gradient Computation

Recall that gradient $\nabla_w \mathcal{L}(w)$ is computed over (iid) data $\mathcal{D} = \{(x_1, y_1), \ldots, (x_N, y_N)\}$.

$$\nabla_w \mathcal{L}(w) = \nabla_w \left[ -\frac{1}{N} \sum_{i=1}^{N} \log p(y_i | x_i, w) \right] \qquad \text{(e.g. logistic regression)}$$

$$= -\frac{1}{N} \sum_{i=1}^{N} \nabla_w \log p(y_i | x_i, w) = -\frac{1}{N} \sum_{i=1}^{N} \nabla_w \mathcal{L}^i(w)$$

### Challenge

- Estimation requires evaluating gradients at *all* $N$ data points
- Fine if $N$ is small, but if $N$ is large? Say millions?
- Can we get a "good enough" gradient from fewer data points? Maybe one!?

## Stochastic Gradient Descent

**Idea:** Compute update for parameter with just a single instance

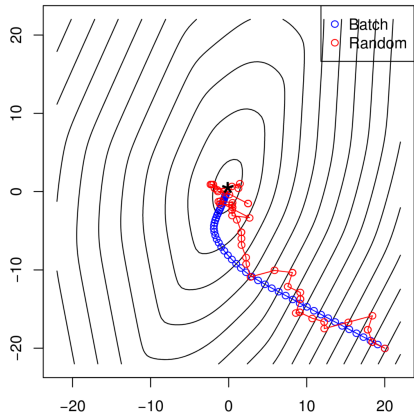$$w \leftarrow w - \eta \cdot \nabla_w \mathcal{L}^i(w)$$

### Indexing

- Choose randomly for $i \in \{1, \ldots, N\}$
- $\mathbb{E}\left[\nabla_w \mathcal{L}^i(w)\right] = \nabla_w \mathcal{L}(w)$
- provides an unbiased estimate of the gradient at each step

### Features

- Cost per iteration independent of $N$
- Potential savings in memory usage
- Can be noisy in practice

# Stochastic Gradient Descent — Example

Example with $N = 10$, $D = 2$ comparing standard versus stochastic gradient descent

**Blue** standard steps $O(N \times D)$

**Red** stochastic steps $O(D)$

## Characteristics

- work well far from optimum
- struggle close to optimum

---

# Stochastic Gradient Descent — Mini-batches

**Idea:** Compute update for parameter with a *few* random instances chosen as $I \subseteq \{1, \ldots, N\}$, such that $|I| = B$, $B \ll N$.

$$w \leftarrow w - \eta \cdot \frac{1}{B} \sum_{i \in I} \nabla_w \mathcal{L}^i(w)$$

Again, we are approximating the full gradient using an unbiased estimate

$$\mathbb{E}\left[\frac{1}{B} \sum_{i \in I} \nabla_w \mathcal{L}^i(w)\right] = \nabla_w \mathcal{L}(w)$$

## Features

- reduces the variance of the gradient estimator by $1/B$
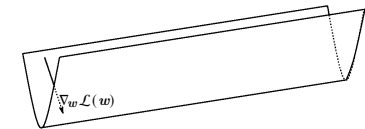- also $B$ times more expensive to compute $O(B \times D)$

---

# Gradient Descent

## Problems

---

# Problems With Gradient Descent

- Setting the step size $\eta$
- Shallow valleys
- Surface curvature
- Local minima



- Gradient descent slows down in shallow valleys
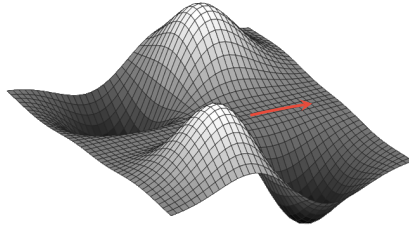- Incorporate *momentum*

$$d \leftarrow \beta d + (1 - \beta)\eta \cdot \nabla_w \mathcal{L}(w)$$

- Have to choose both $\eta$ and $\beta$
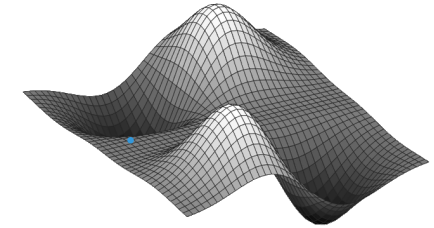
# Problems With Gradient Descent

- Setting the step size $\eta$
- Shallow valleys
- Surface curvature
- Local minima



- Gradient need not point towards optimum!

  **Note:** *locally* steepest direction
- Local curvature is measured by the Hessian: $H = \nabla_w^2 \mathcal{L}(w)$

# Problems With Gradient Descent

- Setting the step size $\eta$
- Shallow valleys
- Surface curvature
- Local minima



- Gradient at *any* minimum is 0!
- *Convex* functions: local minimum $\equiv$ global minimum

  e.g. squared error, logistic regression likelihood, …
- No standard solution

  best to rerun optimiser from different initialisations

# Summary

- Optimisation is a complex field
  - How and why we convert learning problems into optimization problems
  - Gradient Descent / Stochastic Gradient Descent
  - Issues with Gradient Descent
- Many variants providing better stability and convergence
  E.g. momentum, acceleration, averaging, variance reduction, …
- See AdaGrad, Adam, AdaMax, … and many more!