

Basics of Numerical Optimization: Preliminaries

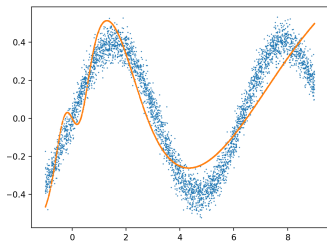
Ju Sun

Computer Science & Engineering

University of Minnesota, Twin Cities

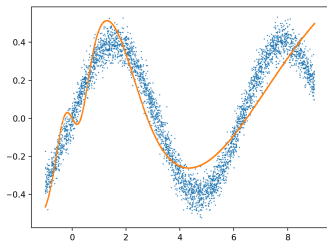
February 11, 2020

Supervised learning as function approximation



- Underlying true function: f_0
- Training data: $\{x_i, y_i\}$ with $y_i \approx f_0(x_i)$
- Choose a family of functions \mathcal{H} , so that
 $\exists f \in \mathcal{H}$ and f and f_0 are close

Supervised learning as function approximation

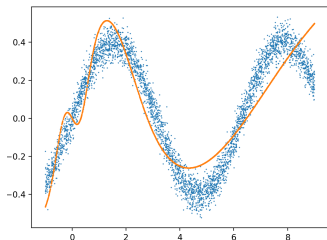


- Underlying true function: f_0
- Training data: $\{x_i, y_i\}$ with $y_i \approx f_0(x_i)$
- Choose a family of functions \mathcal{H} , so that $\exists f \in \mathcal{H}$ and f and f_0 are close
- Find f , i.e., optimization

$$\min_{f \in \mathcal{H}} \sum_i \ell(y_i, f(x_i)) + \Omega(f)$$

- **Approximation capacity: Universal approximation theorems (UAT)**
 \implies replace \mathcal{H} by $\text{DNN}_{\mathbf{W}}$, i.e., a deep neural network with weights \mathbf{W}

Supervised learning as function approximation



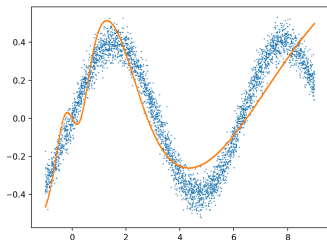
- Underlying true function: f_0
- Training data: $\{x_i, y_i\}$ with $y_i \approx f_0(x_i)$
- Choose a family of functions \mathcal{H} , so that $\exists f \in \mathcal{H}$ and f and f_0 are close
- Find f , i.e., optimization

$$\min_{f \in \mathcal{H}} \sum_i \ell(y_i, f(x_i)) + \Omega(f)$$

- **Approximation capacity: Universal approximation theorems (UAT)**
 \implies replace \mathcal{H} by $\text{DNN}_{\mathbf{W}}$, i.e., a deep neural network with weights \mathbf{W}
- **Optimization:**

$$\min_{\mathbf{W}} \sum_i \ell(y_i, \text{DNN}_{\mathbf{W}}(x_i)) + \Omega(\mathbf{W})$$

Supervised learning as function approximation



- Underlying true function: f_0
- Training data: $\{x_i, y_i\}$ with $y_i \approx f_0(x_i)$
- Choose a family of functions \mathcal{H} , so that $\exists f \in \mathcal{H}$ and f and f_0 are close
- Find f , i.e., optimization

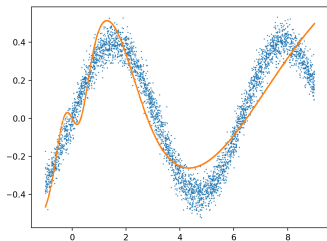
$$\min_{f \in \mathcal{H}} \sum_i \ell(y_i, f(x_i)) + \Omega(f)$$

- **Approximation capacity: Universal approximation theorems (UAT)**
 \implies replace \mathcal{H} by $\text{DNN}_{\mathbf{W}}$, i.e., a deep neural network with weights \mathbf{W}
- **Optimization:**

$$\min_{\mathbf{W}} \sum_i \ell(y_i, \text{DNN}_{\mathbf{W}}(x_i)) + \Omega(\mathbf{W})$$

- **Generalization:** how to avoid over-complicated $\text{DNN}_{\mathbf{W}}$ in view of UAT

Supervised learning as function approximation



- Underlying true function: f_0
- Training data: $\{x_i, y_i\}$ with $y_i \approx f_0(x_i)$
- Choose a family of functions \mathcal{H} , so that $\exists f \in \mathcal{H}$ and f and f_0 are close
- Find f , i.e., optimization

$$\min_{f \in \mathcal{H}} \sum_i \ell(y_i, f(x_i)) + \Omega(f)$$

- **Approximation capacity: Universal approximation theorems (UAT)**
 \implies replace \mathcal{H} by $\text{DNN}_{\mathbf{W}}$, i.e., a deep neural network with weights \mathbf{W}
- **Optimization:**

$$\min_{\mathbf{W}} \sum_i \ell(y_i, \text{DNN}_{\mathbf{W}}(x_i)) + \Omega(\mathbf{W})$$

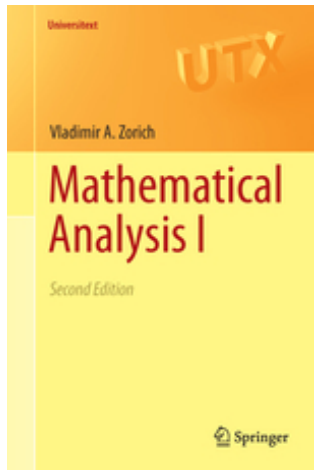
- **Generalization:** how to avoid over-complicated $\text{DNN}_{\mathbf{W}}$ in view of UAT

Now we start to focus on **optimization**.

Elements of multivariate calculus

Optimality conditions of unconstrained optimization

Recommended references



[Munkres, 1997, Zorich, 2015, Coleman, 2012]

Our notation

- scalars: x , vectors: \boldsymbol{x} , matrices: \boldsymbol{X} , tensors: \mathcal{X} , sets: S

Our notation

- scalars: x , vectors: \boldsymbol{x} , matrices: \boldsymbol{X} , tensors: \mathcal{X} , sets: S
- vectors are always **column vectors**, unless stated otherwise

Our notation

- scalars: x , vectors: \mathbf{x} , matrices: \mathbf{X} , tensors: \mathcal{X} , sets: S
- vectors are always **column vectors**, unless stated otherwise
- x_i : i -th element of \mathbf{x} , x_{ij} : (i, j) -th element of \mathbf{X} , \mathbf{x}^i : i -th row of \mathbf{X} as a **row vector**, \mathbf{x}_j : j -th column of \mathbf{X} as a **column vector**

Our notation

- scalars: x , vectors: \mathbf{x} , matrices: \mathbf{X} , tensors: \mathcal{X} , sets: S
- vectors are always **column vectors**, unless stated otherwise
- x_i : i -th element of \mathbf{x} , x_{ij} : (i, j) -th element of \mathbf{X} , \mathbf{x}^i : i -th row of \mathbf{X} as a **row vector**, \mathbf{x}_j : j -th column of \mathbf{X} as a **column vector**
- \mathbb{R} : real numbers, \mathbb{R}_+ : positive reals, \mathbb{R}^n : space of n -dimensional vectors, $\mathbb{R}^{m \times n}$: space of $m \times n$ matrices, $\mathbb{R}^{m \times n \times k}$: space of $m \times n \times k$ tensors, etc

Our notation

- scalars: x , vectors: \mathbf{x} , matrices: \mathbf{X} , tensors: \mathcal{X} , sets: S
- vectors are always **column vectors**, unless stated otherwise
- x_i : i -th element of \mathbf{x} , x_{ij} : (i, j) -th element of \mathbf{X} , \mathbf{x}^i : i -th row of \mathbf{X} as a **row vector**, \mathbf{x}_j : j -th column of \mathbf{X} as a **column vector**
- \mathbb{R} : real numbers, \mathbb{R}_+ : positive reals, \mathbb{R}^n : space of n -dimensional vectors, $\mathbb{R}^{m \times n}$: space of $m \times n$ matrices, $\mathbb{R}^{m \times n \times k}$: space of $m \times n \times k$ tensors, etc
- $[n] \doteq \{1, \dots, n\}$

Differentiability — first order

Consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point x if there exists a matrix $B \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(x + \delta) - f(x) - B\delta}{\|\delta\|_2} \rightarrow \mathbf{0} \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

Differentiability — first order

Consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point x if there exists a matrix $B \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(x + \delta) - f(x) - B\delta}{\|\delta\|_2} \rightarrow \mathbf{0} \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

$$\text{i.e.,} \quad f(x + \delta) = f(x) + B\delta + o(\|\delta\|_2) \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

Differentiability — first order

Consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point x if there exists a matrix $B \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(x + \delta) - f(x) - B\delta}{\|\delta\|_2} \rightarrow \mathbf{0} \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

$$\text{i.e.,} \quad f(x + \delta) = f(x) + B\delta + o(\|\delta\|_2) \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

- B is called the (Fréchet) derivative. When $m = 1$, b^\top (i.e., B^\top) called **gradient**, denoted as $\nabla f(x)$. For general m , also called **Jacobian** matrix, denoted as $J_f(x)$.

Differentiability — first order

Consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point x if there exists a matrix $B \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(x + \delta) - f(x) - B\delta}{\|\delta\|_2} \rightarrow \mathbf{0} \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

$$\text{i.e.,} \quad f(x + \delta) = f(x) + B\delta + o(\|\delta\|_2) \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

- B is called the (Fréchet) derivative. When $m = 1$, b^\top (i.e., B^\top) called **gradient**, denoted as $\nabla f(x)$. For general m , also called **Jacobian** matrix, denoted as $J_f(x)$.
- Calculation: $b_{ij} = \frac{\partial f_i}{\partial x_j}(x)$

Differentiability — first order

Consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point x if there exists a matrix $B \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(x + \delta) - f(x) - B\delta}{\|\delta\|_2} \rightarrow \mathbf{0} \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

$$\text{i.e.,} \quad f(x + \delta) = f(x) + B\delta + o(\|\delta\|_2) \quad \text{as} \quad \delta \rightarrow \mathbf{0}.$$

- B is called the (Fréchet) derivative. When $m = 1$, b^\top (i.e., B^\top) called **gradient**, denoted as $\nabla f(x)$. For general m , also called **Jacobian** matrix, denoted as $J_f(x)$.
- Calculation: $b_{ij} = \frac{\partial f_i}{\partial x_j}(x)$
- **Sufficient condition**: if all partial derivatives exist and are **continuous** at x , then $f(x)$ is differentiable at x .

Calculus rules

Assume $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are differentiable at a point $\mathbf{x} \in \mathbb{R}^n$.

- **linearity:** $\lambda_1 f + \lambda_2 g$ is differentiable at \mathbf{x} and
$$\nabla [\lambda_1 f + \lambda_2 g] (\mathbf{x}) = \lambda_1 \nabla f (\mathbf{x}) + \lambda_2 \nabla g (\mathbf{x})$$
- **product:** assume $m = 1$, fg is differentiable at \mathbf{x} and
$$\nabla [fg] (\mathbf{x}) = f (\mathbf{x}) \nabla g (\mathbf{x}) + g (\mathbf{x}) \nabla f (\mathbf{x})$$
- **quotient:** assume $m = 1$ and $g (\mathbf{x}) \neq 0$, $\frac{f}{g}$ is differentiable at \mathbf{x} and
$$\nabla \left[\frac{f}{g} \right] (\mathbf{x}) = \frac{g(\mathbf{x}) \nabla f(\mathbf{x}) - f(\mathbf{x}) \nabla g(\mathbf{x})}{g^2(\mathbf{x})}$$

Assume $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are differentiable at a point $\mathbf{x} \in \mathbb{R}^n$.

- **linearity:** $\lambda_1 f + \lambda_2 g$ is differentiable at \mathbf{x} and
$$\nabla [\lambda_1 f + \lambda_2 g] (\mathbf{x}) = \lambda_1 \nabla f (\mathbf{x}) + \lambda_2 \nabla g (\mathbf{x})$$
- **product:** assume $m = 1$, fg is differentiable at \mathbf{x} and
$$\nabla [fg] (\mathbf{x}) = f (\mathbf{x}) \nabla g (\mathbf{x}) + g (\mathbf{x}) \nabla f (\mathbf{x})$$
- **quotient:** assume $m = 1$ and $g (\mathbf{x}) \neq 0$, $\frac{f}{g}$ is differentiable at \mathbf{x} and
$$\nabla \left[\frac{f}{g} \right] (\mathbf{x}) = \frac{g (\mathbf{x}) \nabla f (\mathbf{x}) - f (\mathbf{x}) \nabla g (\mathbf{x})}{g^2 (\mathbf{x})}$$
- **Chain rule:** Let $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$, and f is differentiable at \mathbf{x} and $\mathbf{y} = f (\mathbf{x})$ and h is differentiable at \mathbf{y} . Then, $h \circ f : \mathbb{R}^m \rightarrow \mathbb{R}^k$ is differentiable at \mathbf{x} , and

$$\mathbf{J}_{[h \circ f]} (\mathbf{x}) = \mathbf{J}_h (f (\mathbf{x})) \mathbf{J}_f (\mathbf{x}).$$

Calculus rules

Assume $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are differentiable at a point $\mathbf{x} \in \mathbb{R}^n$.

- **linearity:** $\lambda_1 f + \lambda_2 g$ is differentiable at \mathbf{x} and
$$\nabla [\lambda_1 f + \lambda_2 g] (\mathbf{x}) = \lambda_1 \nabla f (\mathbf{x}) + \lambda_2 \nabla g (\mathbf{x})$$
- **product:** assume $m = 1$, $f g$ is differentiable at \mathbf{x} and
$$\nabla [f g] (\mathbf{x}) = f (\mathbf{x}) \nabla g (\mathbf{x}) + g (\mathbf{x}) \nabla f (\mathbf{x})$$
- **quotient:** assume $m = 1$ and $g (\mathbf{x}) \neq 0$, $\frac{f}{g}$ is differentiable at \mathbf{x} and
$$\nabla \left[\frac{f}{g} \right] (\mathbf{x}) = \frac{g (\mathbf{x}) \nabla f (\mathbf{x}) - f (\mathbf{x}) \nabla g (\mathbf{x})}{g^2 (\mathbf{x})}$$
- **Chain rule:** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$, and f is differentiable at \mathbf{x} and $\mathbf{y} = f (\mathbf{x})$ and h is differentiable at \mathbf{y} . Then, $h \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ is differentiable at \mathbf{x} , and

$$\mathbf{J}_{[h \circ f]} (\mathbf{x}) = \mathbf{J}_h (f (\mathbf{x})) \mathbf{J}_f (\mathbf{x}).$$

When $k = 1$,

$$\nabla [h \circ f] (\mathbf{x}) = \mathbf{J}_f^\top (\mathbf{x}) \nabla h (f (\mathbf{x})).$$

Differentiability — second order

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and assume f is 1st-order differentiable in a small ball around \mathbf{x}

- Write $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[\frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$ provided the right side well defined

Differentiability — second order

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and assume f is 1st-order differentiable in a small ball around \mathbf{x}

- Write $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[\frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$ provided the right side well defined
- **Symmetry:** If both $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ and $\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x})$ exist and both are continuous at \mathbf{x} , then **they are equal**.
- **Hessian (matrix):**

$$\nabla^2 f(\mathbf{x}) \doteq \left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i}, \quad (1)$$

where $\left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i} \in \mathbb{R}^{n \times n}$ has its (j, i) -th element as $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$.

Differentiability — second order

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and assume f is 1st-order differentiable in a small ball around \mathbf{x}

- Write $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[\frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$ provided the right side well defined
- **Symmetry:** If both $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ and $\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x})$ exist and both are continuous at \mathbf{x} , then **they are equal**.
- **Hessian (matrix):**

$$\nabla^2 f(\mathbf{x}) \doteq \left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i}, \quad (1)$$

where $\left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i} \in \mathbb{R}^{n \times n}$ has its (j, i) -th element as $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$.

- $\nabla^2 f$ is symmetric.

Differentiability — second order

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and assume f is 1st-order differentiable in a small ball around \mathbf{x}

- Write $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[\frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$ provided the right side well defined
- **Symmetry:** If both $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ and $\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x})$ exist and both are continuous at \mathbf{x} , then **they are equal**.
- **Hessian (matrix):**

$$\nabla^2 f(\mathbf{x}) \doteq \left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i}, \quad (1)$$

where $\left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i} \in \mathbb{R}^{n \times n}$ has its (j, i) -th element as $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$.

- $\nabla^2 f$ is symmetric.
- **Sufficient condition:** if all $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ exist and are **continuous**, f is 2nd-order differentiable at \mathbf{x} (**not converse; we omit the definition due to its technicality**).

Taylor's theorem

Vector version: consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

Taylor's theorem

Vector version: consider $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at x , then

$$f(x + \delta) = f(x) + \langle \nabla f(x), \delta \rangle + o(\|\delta\|_2) \text{ as } \delta \rightarrow 0.$$

- If f is 2nd-order differentiable at x , then

$$f(x + \delta) = f(x) + \langle \nabla f(x), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x) \delta \rangle + o(\|\delta\|_2^2) \text{ as } \delta \rightarrow 0.$$

Taylor's theorem

Vector version: consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

- If f is 2nd-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

Matrix version: consider $f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at \mathbf{X} , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + o(\|\boldsymbol{\Delta}\|_F) \text{ as } \boldsymbol{\Delta} \rightarrow \mathbf{0}.$$

Taylor's theorem

Vector version: consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

- If f is 2nd-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

Matrix version: consider $f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

- If f is 1st-order differentiable at \mathbf{X} , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + o(\|\boldsymbol{\Delta}\|_F) \text{ as } \boldsymbol{\Delta} \rightarrow \mathbf{0}.$$

- If f is 2nd-order differentiable at \mathbf{X} , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + \frac{1}{2} \langle \boldsymbol{\Delta}, \nabla^2 f(\mathbf{X}) \boldsymbol{\Delta} \rangle + o(\|\boldsymbol{\Delta}\|_F^2) \\ \text{as } \boldsymbol{\Delta} \rightarrow \mathbf{0}.$$

Taylor approximation — asymptotic uniqueness

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be k ($k \geq 1$ integer) times differentiable at a point x . If $P(\delta)$ is a k -th order polynomial satisfying $f(x + \delta) - P(\delta) = o(\delta^k)$ as $\delta \rightarrow 0$, then $P(\delta) = P_k(\delta) \doteq f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$.

Taylor approximation — asymptotic uniqueness

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be k ($k \geq 1$ integer) times differentiable at a point x . If $P(\delta)$ is a k -th order polynomial satisfying $f(x + \delta) - P(\delta) = o(\delta^k)$ as $\delta \rightarrow 0$, then $P(\delta) = P_k(\delta) \doteq f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$.

Generalization to the vector version

- Assume $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 1-order differentiable at \mathbf{x} . If $P(\boldsymbol{\delta}) \doteq f(\mathbf{x}) + \langle \mathbf{v}, \boldsymbol{\delta} \rangle$ satisfies that

$$f(\mathbf{x} + \boldsymbol{\delta}) - P(\boldsymbol{\delta}) = o(\|\boldsymbol{\delta}\|_2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0},$$

then $P(\boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$, i.e., the 1st-order Taylor expansion.

Taylor approximation — asymptotic uniqueness

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be k ($k \geq 1$ integer) times differentiable at a point x . If $P(\delta)$ is a k -th order polynomial satisfying $f(x + \delta) - P(\delta) = o(\delta^k)$ as $\delta \rightarrow 0$, then $P(\delta) = P_k(\delta) \doteq f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$.

Generalization to the vector version

- Assume $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 1-order differentiable at x . If $P(\delta) \doteq f(x) + \langle v, \delta \rangle$ satisfies that

$$f(x + \delta) - P(\delta) = o(\|\delta\|_2) \quad \text{as } \delta \rightarrow 0,$$

then $P(\delta) = f(x) + \langle \nabla f(x), \delta \rangle$, i.e., the 1st-order Taylor expansion.

- Assume $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 2-order differentiable at x . If $P(\delta) \doteq f(x) + \langle v, \delta \rangle + \frac{1}{2} \langle \delta, H \delta \rangle$ with H symmetric satisfies that

$$f(x + \delta) - P(\delta) = o(\|\delta\|_2^2) \quad \text{as } \delta \rightarrow 0,$$

then $P(\delta) = f(x) + \langle \nabla f(x), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x) \delta \rangle$, i.e., the 2nd-order Taylor expansion. We can read off ∇f and $\nabla^2 f$ if we know the expansion!

Taylor approximation — asymptotic uniqueness

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be k ($k \geq 1$ integer) times differentiable at a point x . If $P(\delta)$ is a k -th order polynomial satisfying $f(x + \delta) - P(\delta) = o(\delta^k)$ as $\delta \rightarrow 0$, then $P(\delta) = P_k(\delta) \doteq f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$.

Generalization to the vector version

- Assume $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 1-order differentiable at x . If $P(\delta) \doteq f(x) + \langle v, \delta \rangle$ satisfies that

$$f(x + \delta) - P(\delta) = o(\|\delta\|_2) \quad \text{as } \delta \rightarrow 0,$$

then $P(\delta) = f(x) + \langle \nabla f(x), \delta \rangle$, i.e., the 1st-order Taylor expansion.

- Assume $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 2-order differentiable at x . If $P(\delta) \doteq f(x) + \langle v, \delta \rangle + \frac{1}{2} \langle \delta, H \delta \rangle$ with H symmetric satisfies that

$$f(x + \delta) - P(\delta) = o(\|\delta\|_2^2) \quad \text{as } \delta \rightarrow 0,$$

then $P(\delta) = f(x) + \langle \nabla f(x), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x) \delta \rangle$, i.e., the 2nd-order Taylor expansion. **We can read off ∇f and $\nabla^2 f$ if we know the expansion!**

Similarly for the matrix version. See Chap 5 of [Coleman, 2012] for other forms of Taylor theorems and proofs of the asymptotic uniqueness.

Asymptotic uniqueness — why interesting?

Two ways of deriving gradients and Hessians (Recall HW0!)

- (a) Derive the gradient and Hessian of the linear least-squares function $f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. Please include your calculation details.
- (b) Let $\sigma = \frac{1}{1+e^{-x}}$, i.e., the *logistic function*. Derive the gradient of the matrix-variable function $g(\mathbf{W}) = \|\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})\|_2^2$, where σ is applied to the vector $\mathbf{W}\mathbf{x}$ elementwise. This is regression based on a simplified one-neuron network. Please include your calculation details.
- (a) Consider the least-squares objective $f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ again. Recall that for any two vectors \mathbf{a}, \mathbf{b} , $\|\mathbf{a} - \mathbf{b}\|_2^2 = \|\mathbf{a}\|_2^2 - 2\mathbf{a}^\top \mathbf{b} + \|\mathbf{b}\|_2^2$. Now $f(\mathbf{x} + \delta) = \|\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{A}\delta\|_2^2$. Expand this square by the previous formula, and compare it to the 2nd order Taylor expansion by plugging your results from **Problem 1(a)**. Are they equal or not? Why? (Hint: You may find this fact useful: for any two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and any matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\langle \mathbf{u}, \mathbf{M}\mathbf{v} \rangle = \langle \mathbf{M}^\top \mathbf{u}, \mathbf{v} \rangle$. This can be derived from the trace cyclic property above.)
- (b) Consider the one-neuron network regression again: $g(\mathbf{W}) = \|\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})\|_2^2$ with $\sigma = \frac{1}{1+e^{-x}}$, i.e., the *logistic function*. Let's try to work out its 1st order Taylor expansion by direct expansion as follows.
- Show that $\sigma((\mathbf{W} + \Delta)\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x}) + \sigma'(\mathbf{W}\mathbf{x}) \odot (\Delta\mathbf{x}) + o(\|\Delta\|_F)$ when $\Delta \rightarrow \mathbf{0}$. Here, both σ and σ' are applied elementwise, and \odot denotes the elementwise (Hadamard) product.
 - So $\mathbf{y} - \sigma((\mathbf{W} + \Delta)\mathbf{x}) = (\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})) - \sigma'(\mathbf{W}\mathbf{x}) \odot (\Delta\mathbf{x}) - o(\|\Delta\|_F)$ when $\Delta \rightarrow \mathbf{0}$. Substitute this back into the square and use the identity $\|\mathbf{a} + \mathbf{b} + \mathbf{c}\|_2^2 = \|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2 + \|\mathbf{c}\|_2^2 + 2\mathbf{a}^\top \mathbf{b} + 2\mathbf{a}^\top \mathbf{c} + 2\mathbf{b}^\top \mathbf{c}$ to obtain the first-order approximation to $g(\mathbf{W} + \Delta)$. Remember that any terms lower order than $\|\Delta\|_F$ are not interesting and we can always assume Δ as small as needed.
 - Substitute the result from **Problem 1(b)** into the 1st order Taylor expansion formula above and compare it to the result obtained here. Are they equal or not?

Asymptotic uniqueness — why interesting?

Think of neural networks with identity activation functions

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \mathbf{W}_{k-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_F^2$$

Asymptotic uniqueness — why interesting?

Think of neural networks with identity activation functions

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \mathbf{W}_{k-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_F^2$$

How to derive the gradient?

- Scalar chain rule?

Asymptotic uniqueness — why interesting?

Think of neural networks with identity activation functions

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \mathbf{W}_{k-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_F^2$$

How to derive the gradient?

- Scalar chain rule?
- Vector chain rule?

Asymptotic uniqueness — why interesting?

Think of neural networks with identity activation functions

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \mathbf{W}_{k-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_F^2$$

How to derive the gradient?

- Scalar chain rule?
- Vector chain rule?
- First-order Taylor expansion

Asymptotic uniqueness — why interesting?

Think of neural networks with identity activation functions

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \mathbf{W}_{k-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_F^2$$

How to derive the gradient?

- Scalar chain rule?
- Vector chain rule?
- First-order Taylor expansion

Why interesting? See e.g.,

[Kawaguchi, 2016, Lampinen and Ganguli, 2018]

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- **directional derivative:** $D_{\mathbf{v}} f(\mathbf{x}) \doteq \frac{d}{dt} f(\mathbf{x} + t\mathbf{v})$

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- **directional derivative:** $D_{\mathbf{v}}f(\mathbf{x}) \doteq \frac{d}{dt}f(\mathbf{x} + t\mathbf{v})$
- When f is 1-st order differentiable at \mathbf{x} ,

$$D_{\mathbf{v}}f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle .$$

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- **directional derivative:** $D_{\mathbf{v}}f(\mathbf{x}) \doteq \frac{d}{dt}f(\mathbf{x} + t\mathbf{v})$
- When f is 1-st order differentiable at \mathbf{x} ,

$$D_{\mathbf{v}}f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle .$$

- Now $D_{\mathbf{v}}f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, what is $D_{\mathbf{u}}(D_{\mathbf{v}}f)(\mathbf{x})$?

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- **directional derivative:** $D_{\mathbf{v}} f(\mathbf{x}) \doteq \frac{d}{dt} f(\mathbf{x} + t\mathbf{v})$
- When f is 1-st order differentiable at \mathbf{x} ,

$$D_{\mathbf{v}} f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

- Now $D_{\mathbf{v}} f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, what is $D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x})$?

$$D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{v} \rangle.$$

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- **directional derivative:** $D_{\mathbf{v}} f(\mathbf{x}) \doteq \frac{d}{dt} f(\mathbf{x} + t\mathbf{v})$
- When f is 1-st order differentiable at \mathbf{x} ,

$$D_{\mathbf{v}} f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

- Now $D_{\mathbf{v}} f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, what is $D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x})$?

$$D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{v} \rangle.$$

- When $\mathbf{u} = \mathbf{v}$,

$$D_{\mathbf{u}}(D_{\mathbf{u}} f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle = \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{u}).$$

Directional derivatives and curvatures

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

– **directional derivative:** $D_{\mathbf{v}} f(\mathbf{x}) \doteq \frac{d}{dt} f(\mathbf{x} + t\mathbf{v})$

– When f is 1-st order differentiable at \mathbf{x} ,

$$D_{\mathbf{v}} f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

– Now $D_{\mathbf{v}} f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, what is $D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x})$?

$$D_{\mathbf{u}}(D_{\mathbf{v}} f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{v} \rangle.$$

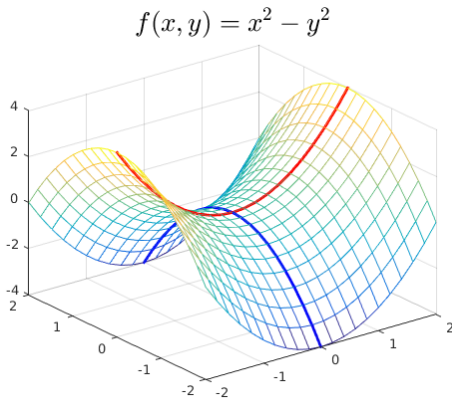
– When $\mathbf{u} = \mathbf{v}$,

$$D_{\mathbf{u}}(D_{\mathbf{u}} f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle = \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{u}).$$

– $\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$ is the **directional curvature** along \mathbf{u} independent of the norm of \mathbf{u}

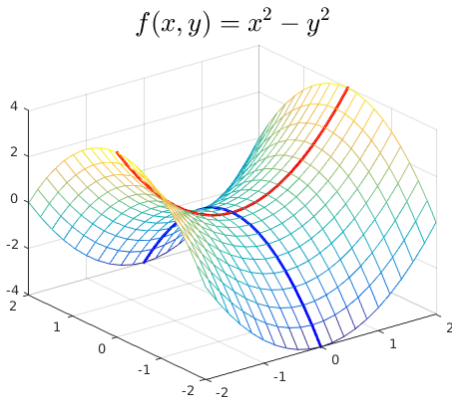
Directional curvature

$\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$ is the **directional curvature** along \mathbf{u} independent of the norm of \mathbf{u}



Directional curvature

$\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$ is the **directional curvature** along \mathbf{u} independent of the norm of \mathbf{u}



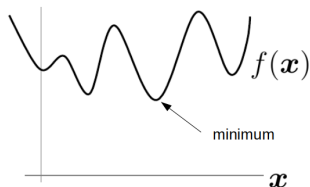
Blue: negative curvature (bending down)

Red: positive curvature (bending up)

Elements of multivariate calculus

Optimality conditions of unconstrained optimization

Optimization problems

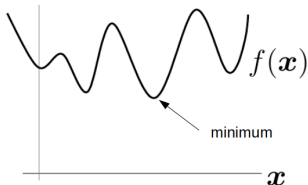


Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler

$$\min_x f(x) \text{ s. t. } x \in C.$$

- x : optimization variables, $f(x)$: objective function, C : constraint (or feasible) set

Optimization problems

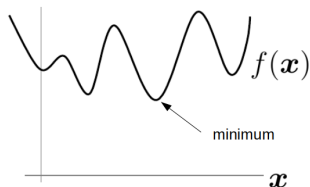


Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler

$$\min_x f(x) \text{ s. t. } x \in C.$$

- x : optimization variables, $f(x)$: objective function, C : constraint (or feasible) set
- C consists of discrete values (e.g., $\{-1, +1\}^n$): discrete optimization; C consists of continuous values (e.g., \mathbb{R}^n , $[0, 1]^n$): **continuous optimization**

Optimization problems

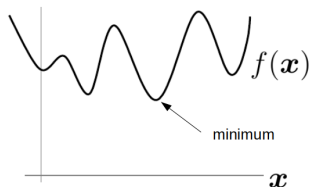


Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler

$$\min_x f(x) \text{ s. t. } x \in C.$$

- x : optimization variables, $f(x)$: objective function, C : constraint (or feasible) set
- C consists of discrete values (e.g., $\{-1, +1\}^n$): discrete optimization; C consists of continuous values (e.g., \mathbb{R}^n , $[0, 1]^n$): **continuous optimization**
- C whole space \mathbb{R}^n : **unconstrained optimization**; C a strict subset of the space: constrained optimization

Optimization problems



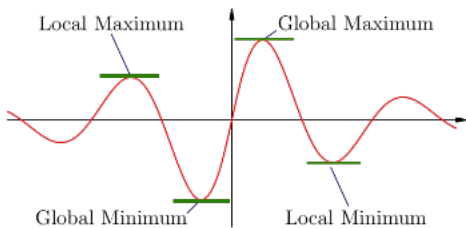
Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler

$$\min_x f(x) \text{ s. t. } x \in C.$$

- x : optimization variables, $f(x)$: objective function, C : constraint (or feasible) set
- C consists of discrete values (e.g., $\{-1, +1\}^n$): discrete optimization; C consists of continuous values (e.g., \mathbb{R}^n , $[0, 1]^n$): **continuous optimization**
- C whole space \mathbb{R}^n : **unconstrained optimization**; C a strict subset of the space: constrained optimization

We focus on **continuous, unconstrained** optimization here.

Global and local mins



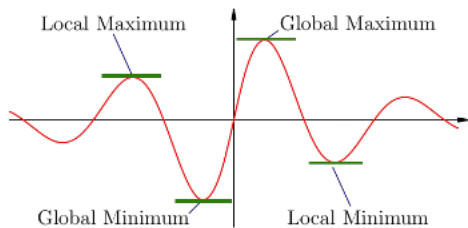
Credit: study.com

Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\min_{x \in \mathbb{R}^n} f(x)$$

- x_0 is a **local minimizer** if: $\exists \varepsilon > 0$, so that $f(x_0) \leq f(x)$ for all x satisfying $\|x - x_0\|_2 < \varepsilon$. The value $f(x_0)$ is called a **local minimum**.

Global and local mins



Credit: study.com

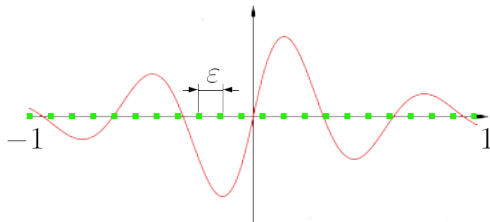
Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\min_{x \in \mathbb{R}^n} f(x)$$

- x_0 is a **local minimizer** if: $\exists \varepsilon > 0$, so that $f(x_0) \leq f(x)$ for all x satisfying $\|x - x_0\|_2 < \varepsilon$. The value $f(x_0)$ is called a **local minimum**.
- x_0 is a **global minimizer** if: $f(x_0) \leq f(x)$ for all $x \in \mathbb{R}^n$. The value is $f(x_0)$ called **the global minimum**.

A naive solution

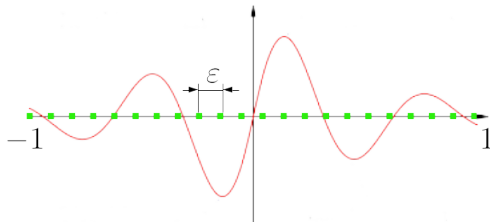
Grid search



- For 1D problem, assume we know the global min lies in $[-1, 1]$

A naive solution

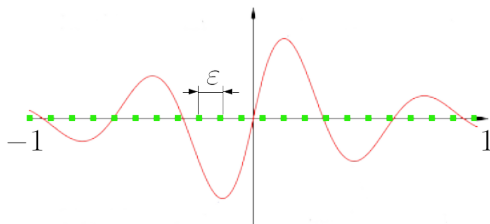
Grid search



- For 1D problem, assume we know the global min lies in $[-1, 1]$
- Take uniformly grid points in $[-1, 1]$ so that any adjacent points are separated by ϵ .

A naive solution

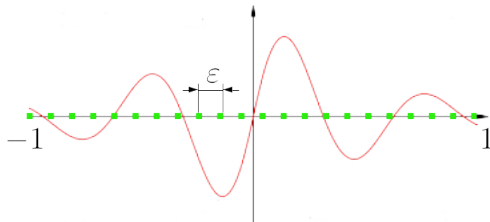
Grid search



- For 1D problem, assume we know the global min lies in $[-1, 1]$
- Take uniformly grid points in $[-1, 1]$ so that any adjacent points are separated by ϵ .
- Need $O(\epsilon^{-1})$ points to get an ϵ -close point to the global min by exhaustive search

A naive solution

Grid search

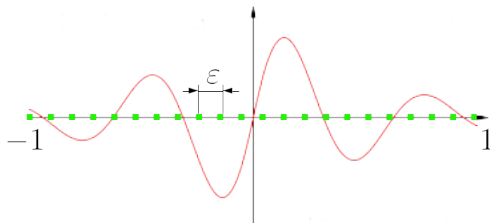


- For 1D problem, assume we know the global min lies in $[-1, 1]$
- Take uniformly grid points in $[-1, 1]$ so that any adjacent points are separated by ϵ .
- Need $O(\epsilon^{-1})$ points to get an ϵ -close point to the global min by exhaustive search

For N -D problems, need $O(\epsilon^{-n})$ computation.

A naive solution

Grid search



- For 1D problem, assume we know the global min lies in $[-1, 1]$
- Take uniformly grid points in $[-1, 1]$ so that any adjacent points are separated by ϵ .
- Need $O(\epsilon^{-1})$ points to get an ϵ -close point to the global min by exhaustive search

For N -D problems, need $O(\epsilon^{-n})$ computation.

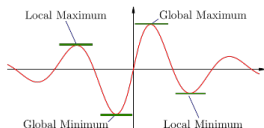
Better characterization of the local/global mins may help avoid this.

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
$$\nabla f(x_0) = \mathbf{0}.$$

First-order optimality condition

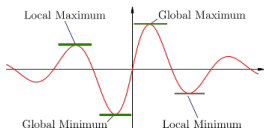
Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
 $\nabla f(x_0) = \mathbf{0}$.



Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer, $\nabla f(x_0) = \mathbf{0}$.

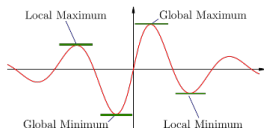


Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$. If x_0 is a local min:

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
 $\nabla f(x_0) = \mathbf{0}$.



Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

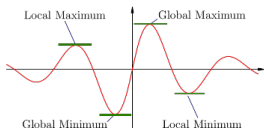
Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$. If x_0 is a local min:

– For all δ sufficiently small,

$$f(x_0 + \delta) - f(x_0) = \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2) \geq 0$$

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
$$\nabla f(x_0) = \mathbf{0}.$$



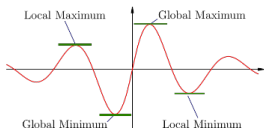
Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$. If x_0 is a local min:

- For all δ sufficiently small,
$$f(x_0 + \delta) - f(x_0) = \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2) \geq 0$$
- For all δ sufficiently small, sign of $\langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$ determined by the sign of $\langle \nabla f(x_0), \delta \rangle$, i.e., $\langle \nabla f(x_0), \delta \rangle \geq 0$.

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
 $\nabla f(x_0) = \mathbf{0}$.



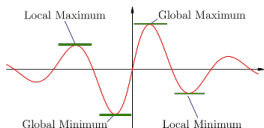
Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$. If x_0 is a local min:

- For all δ sufficiently small,
 $f(x_0 + \delta) - f(x_0) = \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2) \geq 0$
- For all δ sufficiently small, sign of $\langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$ determined by the sign of $\langle \nabla f(x_0), \delta \rangle$, i.e., $\langle \nabla f(x_0), \delta \rangle \geq 0$.
- So for all δ sufficiently small, $\langle \nabla f(x_0), \delta \rangle \geq 0$ and
 $\langle \nabla f(x_0), -\delta \rangle = -\langle \nabla f(x_0), \delta \rangle \geq 0 \implies \langle \nabla f(x_0), \delta \rangle = 0$

First-order optimality condition

Assume f is 1st-order differentiable at x_0 . If x_0 is a local minimizer,
 $\nabla f(x_0) = \mathbf{0}$.



Intuition: ∇f is “rate of change” of function value. If the rate is not zero at x_0 , possible to decrease f along $-\nabla f(x_0)$

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$. If x_0 is a local min:

- For all δ sufficiently small,
 $f(x_0 + \delta) - f(x_0) = \langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2) \geq 0$
- For all δ sufficiently small, sign of $\langle \nabla f(x_0), \delta \rangle + o(\|\delta\|_2)$ determined by the sign of $\langle \nabla f(x_0), \delta \rangle$, i.e., $\langle \nabla f(x_0), \delta \rangle \geq 0$.
- So for all δ sufficiently small, $\langle \nabla f(x_0), \delta \rangle \geq 0$ and
 $\langle \nabla f(x_0), -\delta \rangle = -\langle \nabla f(x_0), \delta \rangle \geq 0 \implies \langle \nabla f(x_0), \delta \rangle = 0$
- So $\nabla f(x_0) = \mathbf{0}$.

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

- If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ (1st-order condition) and $f(x_0 + \delta) = f(x_0) + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

- If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ (1st-order condition) and $f(x_0 + \delta) = f(x_0) + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.
- So $f(x_0 + \delta) - f(x_0) = \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2) \geq 0$ for all δ sufficiently small

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

- If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ (1st-order condition) and $f(x_0 + \delta) = f(x_0) + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.
- So $f(x_0 + \delta) - f(x_0) = \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2) \geq 0$ for all δ sufficiently small
- For all δ sufficiently small, sign of $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$ determined by the sign of $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \implies \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \geq 0$

Second-order optimality condition

Necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$ (i.e., positive semidefinite).

Sufficient condition: Assume $f(x)$ is 2-order differentiable at x_0 . If $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succ \mathbf{0}$ (i.e., positive definite), x_0 is a local min.

Taylor's: $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.

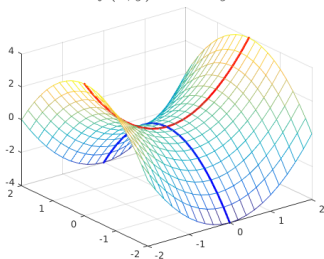
- If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ (1st-order condition) and $f(x_0 + \delta) = f(x_0) + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$.
- So $f(x_0 + \delta) - f(x_0) = \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2) \geq 0$ for all δ sufficiently small
- For all δ sufficiently small, sign of $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$ determined by the sign of $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \implies \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \geq 0$
- So $\nabla^2 f(x_0) \succeq \mathbf{0}$.

What's in between?

2nd order sufficient: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succ \mathbf{0}$

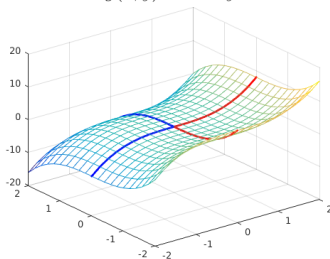
2nd order necessary: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$

$$f(x, y) = x^2 - y^2$$



$$\nabla f = \begin{bmatrix} 2x \\ -2y \end{bmatrix}, \nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

$$g(x, y) = x^3 - y^3$$



$$\nabla g = \begin{bmatrix} 3x^2 \\ -3y^2 \end{bmatrix}, \nabla^2 g = \begin{bmatrix} 6x & 0 \\ 0 & -6y \end{bmatrix}$$

- [Coleman, 2012] Coleman, R. (2012). **Calculus on Normed Vector Spaces**. Springer New York.
- [Kawaguchi, 2016] Kawaguchi, K. (2016). **Deep learning without poor local minima**. *arXiv:1605.07110*.
- [Lampinen and Ganguli, 2018] Lampinen, A. K. and Ganguli, S. (2018). **An analytic theory of generalization dynamics and transfer learning in deep linear networks**. *arXiv:1809.10374*.
- [Munkres, 1997] Munkres, J. R. (1997). **Analysis On Manifolds**. Taylor & Francis Inc.
- [Zorich, 2015] Zorich, V. A. (2015). **Mathematical Analysis I**. Springer Berlin Heidelberg.