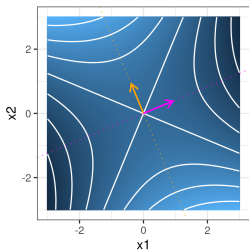# Optimization in Machine Learning

# Mathematical Concepts:
# Quadratic forms II



### Learning goals

- Geometry of quadratic forms
- Spectrum of Hessian

## **PROPERTIES OF QUADRATIC FUNCTIONS**

**Recall**: Quadratic form $q$

- Univariate: $q(x) = ax^2 + bx + c$
- Multivariate: $q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \boldsymbol{b}^T \mathbf{x} + c$

**General observation:** If $q \geq 0$ ($q \leq 0$), $q$ is convex (concave)

**Univariate function:** Second derivative is $q''(x) = 2a$

- $q''(x) \overset{(>)}{\geq} 0$: $q$ (strictly) convex. $q''(x) \overset{(<)}{\leq} 0$: $q$ (strictly) concave.
- High (low) absolute values of $q''(x)$: high (low) curvature

**Multivariate function:** Second derivative is $\mathbf{H} = 2\mathbf{A}$

- Convexity/concavity of $q$ depend on eigenvalues of $\mathbf{H}$
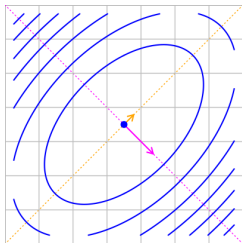- Let us look at an example of the form $q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$

---

# GEOMETRY OF QUADRATIC FUNCTIONS

**Example:** $\mathbf{A} = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \implies \mathbf{H} = 2\mathbf{A} = \begin{pmatrix} 4 & -2 \\ -2 & 4 \end{pmatrix}$

- Since $\mathbf{H}$ symmetric, eigendecomposition $\mathbf{H} = \mathbf{V}\Lambda\mathbf{V}^T$ with

$$\mathbf{V} = \begin{pmatrix} | & | \\ v_{\max} & v_{\min} \\ | & | \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \text{ orthogonal}$$

$$\text{and } \Lambda = \begin{pmatrix} \lambda_{\max} & 0 \\ 0 & \lambda_{\min} \end{pmatrix} = \begin{pmatrix} 6 & 0 \\ 0 & 2 \end{pmatrix}.$$

# GEOMETRY OF QUADRATIC FUNCTIONS

- $\boldsymbol{v}_{max}$ ($\boldsymbol{v}_{min}$) direction of highest (lowest) curvature

  **Proof:** With $\boldsymbol{v} = \mathbf{V}^T\mathbf{x}$:

  $$\mathbf{x}^T\mathbf{H}\mathbf{x} = \mathbf{x}^T\mathbf{V}\Lambda\mathbf{V}^T\mathbf{x} = \boldsymbol{v}^T\Lambda\boldsymbol{v} = \sum_{i=1}^{d} \lambda_i v_i^2 \leq \lambda_{max} \sum_{i=1}^{d} v_i^2 = \lambda_{max}\|\boldsymbol{v}\|^2$$

  Since $\|\boldsymbol{v}\| = \|\mathbf{x}\|$ (**V** orthogonal): $\max_{\|\mathbf{x}\|=1} \mathbf{x}^T\mathbf{H}\mathbf{x} \leq \lambda_{max}$

  Additional: $\boldsymbol{v}_{max}{}^T\mathbf{H}\boldsymbol{v}_{max} = \mathbf{e}_1^T\Lambda\mathbf{e}_1 = \lambda_{max}$

  Analogous: $\min_{\|\mathbf{x}\|=1} \mathbf{x}^T\mathbf{H}\mathbf{x} \geq \lambda_{min}$ and $\boldsymbol{v}_{min}{}^T\mathbf{H}\boldsymbol{v}_{min} = \lambda_{min}$

- $\boldsymbol{v}_{max}, \boldsymbol{v}_{min}$ principal axes of contour ellipses (principal axis theorem)

  **Proof:** With $\boldsymbol{v} = \mathbf{V}^T\mathbf{x}$:

  $$q(\mathbf{x}) = \mathbf{x}^T\mathbf{H}\mathbf{x} + \boldsymbol{b}^T\mathbf{x} + c = \boldsymbol{v}^T\Lambda\boldsymbol{v} + \boldsymbol{b}^T V\boldsymbol{v} + c =: \tilde{q}(\boldsymbol{v})$$

  Now:

  $$q(\boldsymbol{v}_j) = \mathbf{e}_j^T\Lambda\mathbf{e}_j + \boldsymbol{b}^T V\mathbf{e}_j + c = \tilde{q}(\mathbf{e}_j)$$

  Especially: $q(\boldsymbol{v}_{max}) = \lambda_{max} = \tilde{q}(\mathbf{e}_1)$ and $q(\boldsymbol{v}_{min}) = \lambda_{min} = \tilde{q}(\mathbf{e}_d)$

---

## GEOMETRY OF QUADRATIC FUNCTIONS

Recall: **Second order condition for optimality** is **sufficient**.

We skipped the **proof** at first, but can now catch up on it.

If $H(\mathbf{x}^*) \succ 0$ at stationary point $\mathbf{x}^*$, then $\mathbf{x}^*$ is local minimum ($\prec$ for maximum).

**Proof:** Let $\lambda_{\min} > 0$ denote the smallest eigenvalue of $H(\mathbf{x}^*)$. Then:

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \underbrace{\nabla f(\mathbf{x}^*)}_{=0}{}^T(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}\underbrace{(\mathbf{x} - \mathbf{x}^*)^T H(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)}_{\geq \lambda_{\min}\|\mathbf{x}-\mathbf{x}^*\|^2 \text{ (see above)}} + \underbrace{R_2(\mathbf{x}, \mathbf{x}^*)}_{=o(\|\mathbf{x}-\mathbf{x}^*\|^2)} \ .$$
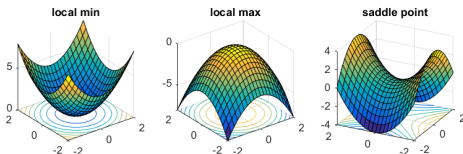
Choose $\epsilon > 0$ s.t. $|R_2(\mathbf{x}, \mathbf{x}^*)| < \frac{1}{2}\lambda_{\min}\|\mathbf{x} - \mathbf{x}^*\|^2$ for each $\mathbf{x}$ with $\|\mathbf{x} - \mathbf{x}^*\| < \epsilon$. Then:

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \underbrace{\frac{1}{2}\lambda_{\min}\|\mathbf{x} - \mathbf{x}^*\|^2 + R_2(\mathbf{x}, \mathbf{x}^*)}_{>0} > f(\mathbf{x}^*) \quad \text{for each } \mathbf{x} \text{ with } \|\mathbf{x} - \mathbf{x}^*\| < \epsilon.$$

# GEOMETRY OF QUADRATIC FUNCTIONS

If spectrum of **A** is known, also that of **H** = 2**A** is known.

- If **all** eigenvalues of $\mathbf{H} \overset{(>)}{\geq} 0$ ($\Leftrightarrow \mathbf{H} \overset{(\succ)}{\succeq} 0$):
    - $q$ (strictly) convex,
    - there is a (unique) global minimum.
- If **all** eigenvalues of $\mathbf{H} \overset{(<)}{\leq} 0$ ($\Leftrightarrow \mathbf{H} \overset{(\prec)}{\preceq} 0$):
    - $q$ (strictly) concave,
    - there is a (unique) global maximum.
- If **H** has both positive and negative eigenvalues ($\Leftrightarrow$ **H** indefinite):
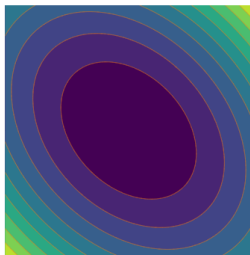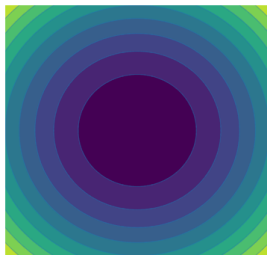    - $q$ neither convex nor concave,
    - there is a saddle point.

# CONDITION AND CURVATURE

Condition of $\mathbf{H} = 2\mathbf{A}$ is given by $\kappa(\mathbf{H}) = \kappa(\mathbf{A}) = |\lambda_{max}|/|\lambda_{min}|$.

**High condition** means:

- $|\lambda_{max}| \gg |\lambda_{min}|$
- Curvature along $\mathbf{v}_{max} \gg$ curvature along $\mathbf{v}_{min}$
- **Problem** for optimization algorithms like **gradient descent** (later)
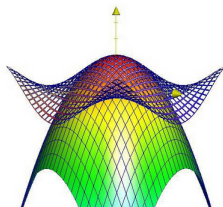


**Left:** Excellent condition. **Middle:** Good condition. **Right:** Bad condition.

# APPROXIMATION OF SMOOTH FUNCTIONS

Any function $f \in \mathcal{C}^2$ can be locally approximated by a quadratic function via second order Taylor approximation:

$$f(\mathbf{x}) \approx f(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}})^T (\mathbf{x} - \tilde{\mathbf{x}}) + \frac{1}{2}(\mathbf{x} - \tilde{\mathbf{x}})^T \nabla^2 f(\tilde{\mathbf{x}})(\mathbf{x} - \tilde{\mathbf{x}})$$



*f* and its second order approximation is shown by the dark and bright grid, respectively. (Source: `daniloroccatano.blog`)

$\implies$ Hessians provide information about **local** geometry of a function.