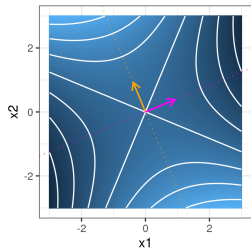


# Optimization in Machine Learning

## Mathematical Concepts

### Quadratic functions II

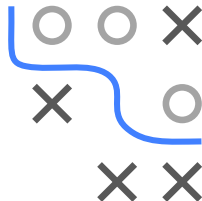


#### Learning goals

- Geometry of quadratic functions
- Spectrum of Hessian

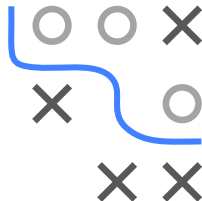
# PROPERTIES OF QUADRATIC FUNCTIONS

- $q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$
- Under symmetry:  $\mathbf{H} = 2\mathbf{A}$
- Convexity/concavity of  $q$  depend on eigenvalues of  $\mathbf{H}$





# SPECTRUM AND CURVATURE

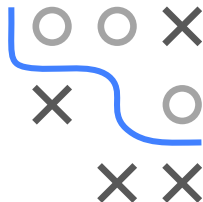


- $\mathbf{v}_{\max}$  direction of highest curvature, with curvature value  $\lambda_{\max}$

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

- Since  $\|\mathbf{v}\| = \|\mathbf{x}\|$  ( $\mathbf{V}$  orthogonal):  $\max_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{H} \mathbf{v} \leq \lambda_{\max}$
- For  $\mathbf{v}_{\max}$  we obtain this upper bound:  $\mathbf{v}_{\max}^T \mathbf{H} \mathbf{v}_{\max} = \mathbf{e}_1^T \mathbf{\Lambda} \mathbf{e}_1 = \lambda_{\max}$
- Analogously,  $\mathbf{v}_{\min}$  direction of lowest curvature, with curvature value  $\lambda_{\min}$
- Contour lines of any quadratic function are ellipses

# SECOND ORDER CONDITION



- Recall: Second order condition for optimality is sufficient
- If  $H(\mathbf{x}^*) \succ 0$  at stationary point  $\mathbf{x}^*$ , then  $\mathbf{x}^*$  local minimum ( $\prec$  for maximum)

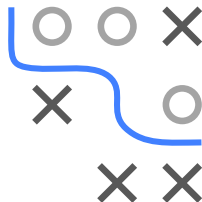
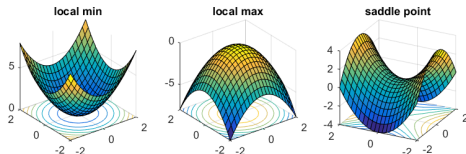
$$f(\mathbf{x}) = f(\mathbf{x}^*) + \underbrace{\nabla f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)}_{=0} + \underbrace{\frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T H(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)}_{\geq \lambda_{\min} \|\mathbf{x} - \mathbf{x}^*\|^2} + \underbrace{R_2(\mathbf{x}, \mathbf{x}^*)}_{=o(\|\mathbf{x} - \mathbf{x}^*\|^2)}$$

- Choose  $\epsilon > 0$  s.t.  $|R_2(\mathbf{x}, \mathbf{x}^*)| < \frac{1}{2} \lambda_{\min} \|\mathbf{x} - \mathbf{x}^*\|^2$  for  $\mathbf{x} \neq \mathbf{x}^*$ ,  $\|\mathbf{x} - \mathbf{x}^*\| < \epsilon$

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \underbrace{\frac{1}{2} \lambda_{\min} \|\mathbf{x} - \mathbf{x}^*\|^2}_{>0} + R_2(\mathbf{x}, \mathbf{x}^*) > f(\mathbf{x}^*)$$

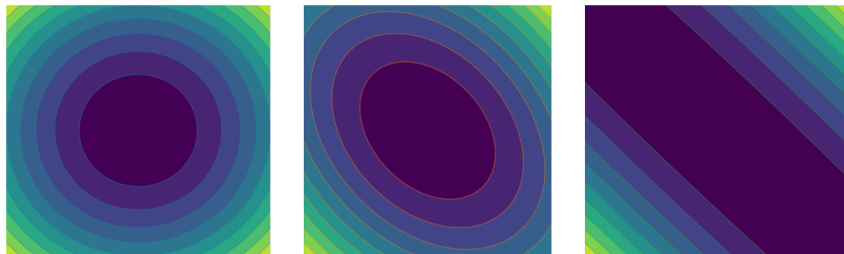
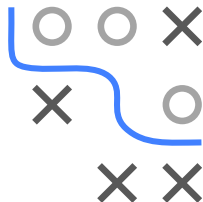
# EIGENVALUES AND SHAPE

- If spectrum of  $\mathbf{A}$  is known, also that of  $\mathbf{H} = 2\mathbf{A}$  is known
- If all eigenvalues of  $\mathbf{H} \stackrel{(>)}{\geq} 0 \Leftrightarrow \mathbf{H} \stackrel{(>)}{\succ} 0$ :
  - $q$  (strictly) convex
  - (Unique) global minimum
- If all eigenvalues of  $\mathbf{H} \stackrel{(<)}{\leq} 0 \Leftrightarrow \mathbf{H} \stackrel{(<)}{\preceq} 0$ :
  - $q$  (strictly) concave
  - (Unique) global maximum
- If  $\mathbf{H}$  has both positive and negative eigenvalues ( $\Leftrightarrow \mathbf{H}$  indefinite):
  - $q$  neither convex nor concave
  - there is a saddle point



# CONDITION AND CURVATURE

- $\kappa(\mathbf{H}) = \kappa(\mathbf{A}) = |\lambda_{\max}|/|\lambda_{\min}|$
- High condition means
  - $|\lambda_{\max}| \gg |\lambda_{\min}|$
  - Curvature along  $\mathbf{v}_{\max} \gg$  along  $\mathbf{v}_{\min}$
  - Problem for algorithms like gradient descent

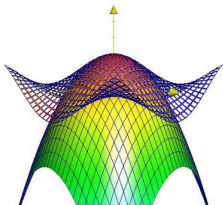


Left: Excellent condition. Middle: Good condition. Right: Bad condition.

# APPROXIMATION OF SMOOTH FUNCTIONS

- Any  $f \in \mathcal{C}^2$  can be locally approximated by quadratic function (second order Taylor)

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



$f$  and second order approximation: dark vs bright grid. (Source: daniloroccatano.blog)

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

