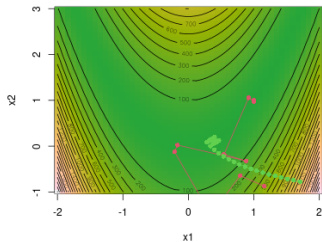# Optimization in Machine Learning

## Second order methods
## Quasi-Newton



**Learning goals**

- Newton-Raphson vs. Quasi-Newton
- SR1
- BFGS

# QUASI-NEWTON: IDEA

Start point of **QN method** is (as with NR) a Taylor approximation of the gradient, except that H is replaced by a **pd** matrix $\boldsymbol{A}^{[t]}$:
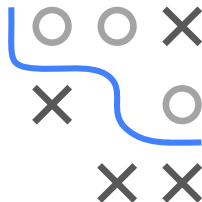
$$\nabla f(\mathbf{x}) \approx \nabla f(\mathbf{x}^{[t]}) + \nabla^2 f(\mathbf{x}^{[t]})(\mathbf{x} - \mathbf{x}^{[t]}) = \mathbf{0} \qquad \text{NR}$$
$$\nabla f(\mathbf{x}) \approx \nabla f(\mathbf{x}^{[t]}) + \boldsymbol{A}^{[t]} \qquad (\mathbf{x} - \mathbf{x}^{[t]}) = \mathbf{0} \qquad \text{QN}$$

The update direction:

$$\boldsymbol{d}^{[t]} = -\nabla^2 f(\mathbf{x}^{[t]})^{-1} \nabla f(\mathbf{x}^{[t]}) \qquad \text{NR}$$
$$\boldsymbol{d}^{[t]} = -(\boldsymbol{A}^{[t]})^{-1} \quad \nabla f(\mathbf{x}^{[t]}) \qquad \text{QN}$$
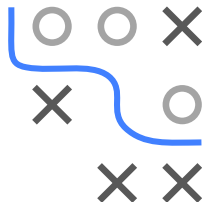
# SYMMETRIC RANK 1 UPDATE (SR1)

Simplest approach: symmetric rank 1 updates (**SR1**) of form

$$\boldsymbol{A}^{[t+1]} \leftarrow \boldsymbol{A}^{[t]} + \boldsymbol{B}^{[t]} = \boldsymbol{A}^{[t]} + \beta \boldsymbol{u}^{[t]} (\boldsymbol{u}^{[t]})^{\top}$$

with appropriate vector $\boldsymbol{u}^{[t]} \in \mathbb{R}^n$, $\beta \in \mathbb{R}$.

## BFGS ALGORITHM

Instead of Rank 1 updates, the **BFGS** procedure (published simultaneously in 1970 by Broyden, Fletcher, Goldfarb and Shanno) uses rank 2 modifications of the form

$$\boldsymbol{A}^{[t]} + \beta_1 \boldsymbol{u}^{[t]} (\boldsymbol{u}^{[t]})^\top + \beta_2 \boldsymbol{v}^{[t]} (\boldsymbol{v}^{[t]})^\top$$

with $\boldsymbol{s}^{[t]} := \boldsymbol{x}^{[t+1]} - \boldsymbol{x}^{[t]}$

- $\boldsymbol{u}^{[t]} = \nabla f(\boldsymbol{x}^{[t+1]}) - \nabla f(\boldsymbol{x}^{[t]})$
- $\boldsymbol{v}^{[t]} = \boldsymbol{A}^{[t]} \boldsymbol{s}^{[t]}$
- $\beta_1 = \frac{1}{(\boldsymbol{u}^{[t]})^\top (\boldsymbol{s}^{[t]})}$
- $\beta_2 = -\frac{1}{(\boldsymbol{s}^{[t]})^\top \boldsymbol{A}^{[t]} \boldsymbol{s}^{[t]}}$

The resulting matrices $\boldsymbol{A}^{[t]}$ are positive definite and the corresponding quasi-newton update directions $\boldsymbol{d}^{[t]}$ are actual descent directions.