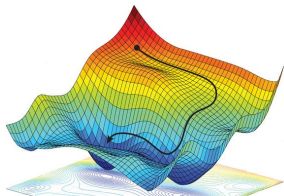


Introduction to Machine Learning

Pseudo-residuals and Gradient Descent



Learning goals

- Learn what pseudo-residuals are
- Understand the relationship between pseudo-residuals and gradient descent

PSEUDO-RESIDUALS

- We further define **pseudo-residuals** as the negative first derivatives of loss functions w.r.t. $f(\mathbf{x})$

$$\tilde{r} := -\frac{\partial L(y, f(\mathbf{x}))}{\partial f(\mathbf{x})}.$$

- We will gain more intuition about the principle of pseudo-residuals in a later chapter.

GD IN ML AND PSEUDO-RESIDUALS

By using the chain rule we see that

$$\begin{aligned}\nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta) &= \sum_{i=1}^n \underbrace{\frac{\partial L(y^{(i)}, f)}{\partial f} \bigg|_{f=f(\mathbf{x}^{(i)} | \theta)}}_{=-\tilde{r}^{(i)}} \cdot \nabla_{\theta} f(\mathbf{x}^{(i)} | \theta) \\ &= - \sum_{i=1}^n \tilde{r}^{(i)} \cdot \nabla_{\theta} f(\mathbf{x}^{(i)} | \theta)\end{aligned}$$

For risk minimization, the update rule for the parameter θ is

$$\begin{aligned}\theta^{[t+1]} &\leftarrow \theta^{[t]} - \alpha^{[t]} \sum_{i=1}^n \nabla_{\theta} L(y^{(i)}, f(\mathbf{x}^{(i)} | \theta)) \bigg|_{\theta=\theta^{[t]}} \\ \theta^{[t+1]} &\leftarrow \theta^{[t]} + \alpha^{[t]} \sum_{i=1}^n \tilde{r}^{(i)} \cdot \nabla_{\theta} f(\mathbf{x}^{(i)} | \theta) \bigg|_{\theta=\theta^{[t]}}\end{aligned}$$

$\alpha^{[t]} \in [0, 1]$ is called “learning rate” in this context.