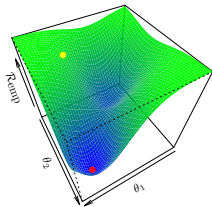


Introduction to Machine Learning

Introduction: Optimization

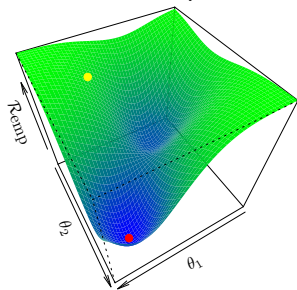


Learning goals

- Understand the difference between local and global minima
- Know the least squares estimator
- Understand the idea of gradient descent

INTRODUCTION

- As we have seen, we can identify models f with their parameters $\theta \in \Theta$ regarding the respective parametrization.
- Hence we can express the associated empirical risk of the model as a function of these parameters.
- Therefore, when we try to find the best model, we actually traverse on the error surface from a starting point (yellow) with the goal of finding the point with the lowest empirical risk (red).



INTRODUCTION

Formally, this means that we find the best model \hat{f} parametrized by parameters $\hat{\theta} \in \Theta$ regarding an empirical risk \mathcal{R}_{emp} by **minimizing** $\mathcal{R}_{\text{emp}}(\theta)$ with respect to θ , i.e.,

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \mathcal{R}_{\text{emp}}(\theta).$$

For such a **(global) minimum** $\hat{\theta}$ it obviously holds that

$$\forall \theta \in \Theta : \quad \mathcal{R}_{\text{emp}}(\hat{\theta}) \leq \mathcal{R}_{\text{emp}}(\theta).$$

However, this does not imply that $\hat{\theta}$ is unique by any means.

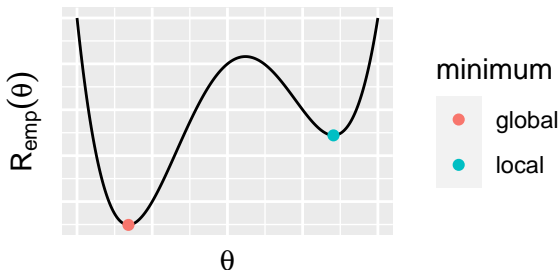
Which kind of technique we can use to solve the minimization problem strongly depends on the feature space. In this chapter we will focus on purely numeric features.

CONTINUOUS \mathcal{R}_{emp}

If the empirical risk \mathcal{R}_{emp} is continuous in θ we can define a **local minimum** $\hat{\theta}$, such that

$$\exists \epsilon > 0 \forall \theta \in \left\{ \bar{\theta} \in \Theta \mid \left\| \hat{\theta} - \bar{\theta} \right\| < \epsilon \right\} : \mathcal{R}_{\text{emp}}(\hat{\theta}) \leq \mathcal{R}_{\text{emp}}(\theta).$$

Clearly every global minimum is also a local minimum (if it exists).
In general finding a local minimum is easier than finding a global minimum.

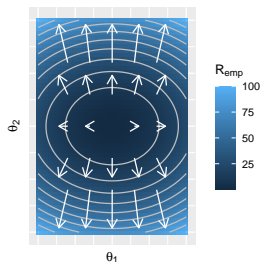


CONTINUOUSLY DIFFERENTIABLE \mathcal{R}_{emp}

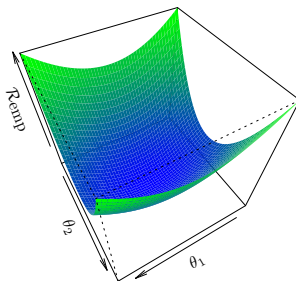
If the empirical risk \mathcal{R}_{emp} is continuously differentiable in θ then a **sufficient condition** for $\hat{\theta}$ to be a local minimum is that the gradient

$$\frac{\partial}{\partial \theta} \mathcal{R}_{\text{emp}}(\hat{\theta}) = 0$$

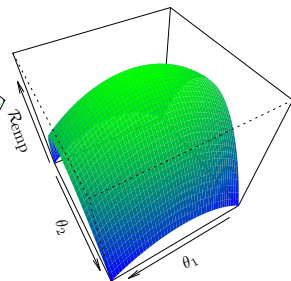
and the Hessian $\frac{\partial^2}{\partial \theta^2} \mathcal{R}_{\text{emp}}(\hat{\theta})$ is positive definite. Which makes sense, since, while the gradient can be thought of as the local direction and rate of fastest increase, the Hessian measures the local curvature of \mathcal{R}_{emp} .



$$0.1 \cdot \frac{\partial}{\partial \theta} \mathcal{R}_{\text{emp}}(\hat{\theta})$$



const. pos. def. Hessian



const. neg. def. Hessian

LEAST SQUARES ESTIMATOR

Now, for given features $\mathbf{X} \in \mathbb{R}^{n \times p}$ and target $\mathbf{y} \in \mathbb{R}^n$, we want to find the best linear model regarding the squared error loss, i.e.,

$$\mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \sum_{i=1}^n (\boldsymbol{\theta}^\top \mathbf{x}^{(i)} - y^{(i)})^2.$$

With the sufficient condition for continuously differentiable functions it can be shown that the **least squares estimator**

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

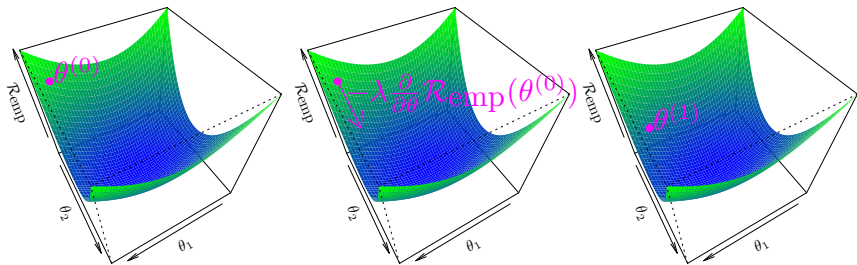
is a local minimum of \mathcal{R}_{emp} . Since, here, \mathcal{R}_{emp} is a convex function it follows that there is only one minimum. Hence $\hat{\boldsymbol{\theta}}$ is the global minimum.

Note: Often such an analytical solution to our respective minimization problem does not exist. Therefore we need numerical methods which enable us to find an approximate solution.

GRADIENT DESCENT

The simple idea of **gradient descent** (GD) is to follow iteratively from the i -th solution candidate $\theta^{(i)}$ in the direction of the negative gradient, i.e., the direction of the steepest descent, with a learning rate λ to the $(i + 1)$ -th solution candidate $\theta^{(i+1)}$, s.t.

$$\theta^{(i+1)} = \theta^{(i)} - \lambda \frac{\partial}{\partial \theta} \mathcal{R}_{\text{emp}}(\theta^{(i)}).$$



FURTHER TOPICS

- There exist many improvements of the GD method, e.g., we could also optimize the learning rate λ .
- GD is a so-called first-order method. Second-order methods use the Hessian (which must therefore exist) to refine the search direction.
- If the gradient of GD is not derived from the empirical risk of the whole data set, but instead from a randomly selected subset of it, we call the respective method **stochastic gradient descent** (SGD). For high-dimensional problems this can lead to a higher computational efficiency.
- Often it is desirable to not allow arbitrarily large $\|\hat{\theta}\|$, since this could result, among other things, in numerical instability of the method. This procedure is called **regularization**.