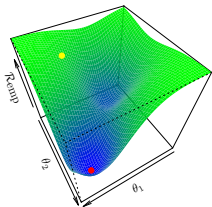


Introduction to Machine Learning

ML-Basics

Optimization

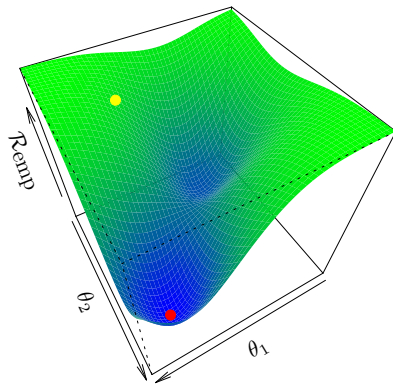


Learning goals

- Understand how the risk function is optimized to learn the optimal parameters of a model
- Understand the idea of gradient descent as a basic risk optimizer

LEARNING AS PARAMETER OPTIMIZATION

- Operationalize search for model f that matches training data best by looking for parametrization $\theta \in \Theta$ with lowest risk $\mathcal{R}_{\text{emp}}(\theta)$
- Traverse error surface downwards; often local search from some start point to minimum (hopefully)



LEARNING AS PARAMETER OPTIMIZATION

ERM optimization problem:

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

For **(global) minimum** $\hat{\theta}$:

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

Does not imply that $\hat{\theta}$ is unique

- Best numerical optimizer depends on problem structure
- Continuous params? Uni-modal $\mathcal{R}_{\text{emp}}(\theta)$?
- Numerical optimization not our focus here, now

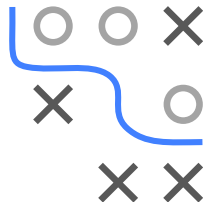
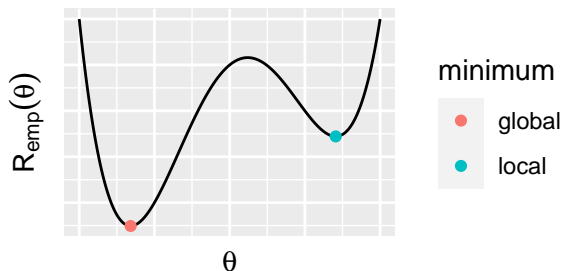


LOCAL MINIMA

- Definition of **local minimum** $\hat{\theta}$:

$$\exists \epsilon > 0 \forall \theta \text{ with } \|\hat{\theta} - \theta\| < \epsilon : \mathcal{R}_{\text{emp}}(\hat{\theta}) \leq \mathcal{R}_{\text{emp}}(\theta)$$

- Clearly every global minimum is also a local minimum
- Finding local minimum is easier than global one



LOCAL MINIMA AND STATIONARY POINTS

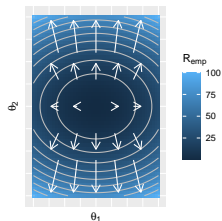
If \mathcal{R}_{emp} continuously differentiable, **sufficient condition** for local minimum:
 $\hat{\theta}$ is **stationary**, so 0 gradient, so no local improvement possible:

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

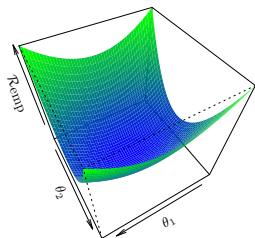
and Hessian at $\hat{\theta}$ is positive definite.

Neg. gradient points into direction of fastest local decrease;

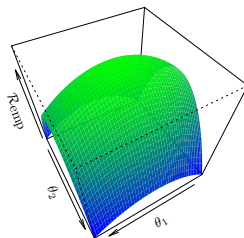
Hessian measures local curvature.



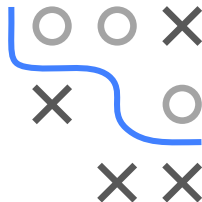
$$\frac{d\mathcal{R}_{\text{emp}}}{d\theta}(\theta)$$



const. pos. def. Hessian



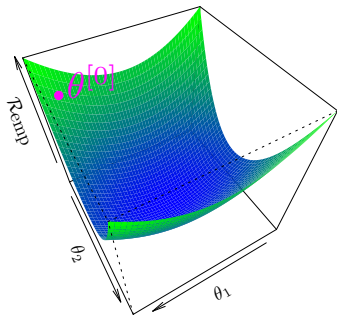
const. neg. def. Hessian



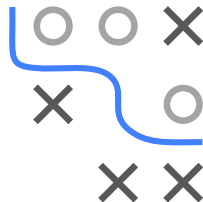
GRADIENT DESCENT

- Iteratively improve current candidate $\theta^{[t]}$
- Move in direction of neg. gradient, so direction of steepest descent
- Use step size / learning rate α

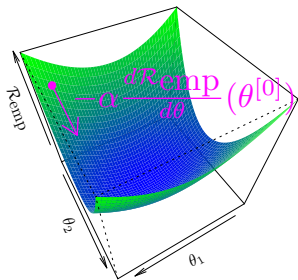
$$\theta^{[t+1]} = \theta^{[t]} - \alpha \frac{d\mathcal{R}_{\text{emp}}}{d\theta}(\theta^{[t]})$$



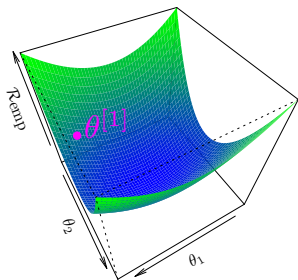
Random start $\theta^{[0]}$ with
 $\mathcal{R}_{\text{emp}}(\theta^{[0]}) = 76.25$.



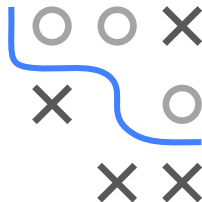
GRADIENT DESCENT - EXAMPLE



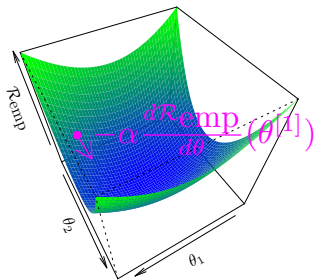
Direction of the neg. gradient at $\theta^{[0]}$



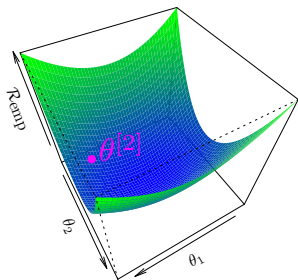
Arrive at $\theta^{[1]}$ with $\mathcal{R}_{\text{emp}}(\theta^{[1]}) \approx 42.73$
We improved: $\mathcal{R}_{\text{emp}}(\theta^{[1]}) < \mathcal{R}_{\text{emp}}(\theta^{[0]})$



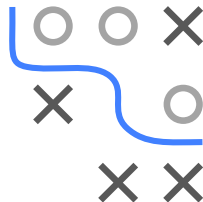
GRADIENT DESCENT - EXAMPLE



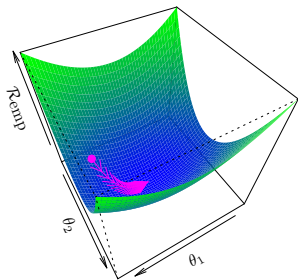
Now iterate, do the same at $\theta^{[1]}$



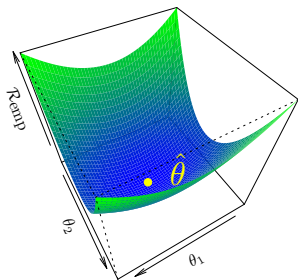
Now $\theta^{[2]}$ has risk $\mathcal{R}_{\text{emp}}(\theta^{[2]}) \approx 25.08$



GRADIENT DESCENT - EXAMPLE



We iterate this until some form of convergence or termination



We arrive close to a stationary $\hat{\theta}$ which is hopefully at least a local minimum



