

ML - How to evaluate Models in terms of finding optimal parameter vector / fitted model

LOSS | Loss function $L: \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}$, $(y, f(x)) \mapsto L(y, f(x))$

Loss function quantifies point-wise(!) the quality of prediction for a single feature vector x by assigning loss/error to it \hat{y} deviation from model predict.

RISK | Theoretical Risk: Expected loss of model $R(f) = \mathbb{E}_{x,y} [L(y, f(x))] = \int L(y, f(x)) dP_{x,y}$

replace integral with finite sum

Empirical Risk: $R_{\text{emp}}(f) := \sum_i L(y^{(i)}, f(x^{(i)}))$ This "quality score" encodes how well f fits training data by measure on i.i.d. training data (θ) parametrised $R_{\text{emp}}: \mathbb{R}^d \rightarrow \mathbb{R}$ summing up all (point) losses
average loss $\bar{R}_{\text{emp}}(\theta) := \frac{1}{n} \dots$ leads to same optimization, average loss is easier to interpret

Empirical Risk Minimization: $\hat{\theta} = \arg \min_{\theta \in \Theta} R_{\text{emp}}(\theta)$ $\hat{\theta}$ not always unique, multiple θ -vectors can have same quality

$$\hat{f} = \arg \min_{f \in \mathcal{F}} R_{\text{emp}}(f)$$

Stochastic gradient descent

Usually we find $\hat{\theta}$, the optimal parameter vector, by numerical optimization methods (Combinatorial optimization, Analytical (OLS normal equation))

(Geometrically) Find minimum in multidimensional risk surface / loss landscape But global minimum does not always exist!

Clearly every global minimum is also a local minimum, which are much easier to find (sometimes those need to suffice)

$$\exists \varepsilon > 0 \forall \theta \text{ with } \|\theta - \hat{\theta}\| \leq \varepsilon : R_{\text{emp}}(\hat{\theta}) \leq R_{\text{emp}}(\theta)$$

→ no local improvement if possible inside very small ε -ball

Necessary condition for local minimum: $\hat{\theta}$ is stationary multidim. "Point": Gradient $\nabla R_{\text{emp}}(\hat{\theta}) = 0$

For continuously differentiable R_{emp} this condition is sufficient if also Hessian Matrix $H(\theta)$ is positive definite

→ positive eigenspectrum

We know: The gradient points into direction of fastest local increase / deepest ascent → quadratic form $x^T H x > 0$

→ negative gradient points into direction of fastest local decrease / steepest descent

GRADIENT DESCENT $\theta^{[t+1]} = \theta^{[t]} - \alpha \nabla R_{\text{emp}}(\theta^{[t]})$ Iterative optimization process for R_{emp} continuous in θ

(Gradient update focus on influential components of θ w.r.t. R_{emp}) \Rightarrow We proceed more quickly along important dims of θ ! Hence up until we reach local minimum

Move from current candidate $\theta^{[t]}$ in direction of negative gradient $\hat{=}$ direction of steepest descent by some defined convergence criterion

When gradient $\rightarrow 0$ we can observe steps becoming super tiny

α is the step-size / learning rate "How fast do we update parameter vector?"

Configuring the step-size α is crucial! We only follow quickest LOCAL decrease, so setting α too large can result in overshooting into unwanted risk increase
Small tiny steps basically guarantee going down, but this is often too slow. We usually strive for something in the middle

We can use constant α , a fixed α -decrease, or dynamic changes of the learning rate

\rightarrow additional computing memory required

GD is a first order method. Second order methods (e.g. Newton-Raphson) use the Hessian to refine search direction for faster convergence

\rightarrow to potentially escape saddle points

There are improvements for GD such as momentum and ADAM. Implementing those may enable us to smartly control the learning rate α , to mimic second order behavior without computing the Hessian etc. Then of course the optimization becomes much more complicated.

Stochastic Gradient Descent (SGD): Gradient is not computed on complete data, but instead on small random batches

For large-scale problems this is usually more efficient