

ML - How to evaluate Models

LOSS] Loss function $L: \mathbb{R}^d \rightarrow \mathbb{R}$, $(y_i, f(x_i)) \mapsto L(y_i, f(x_i))$

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

Observation

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

Empirical Risk: $R_{\text{emp}}(f) := \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$ This "quality score" encodes how well f fits/trains data by measure on i.i.d. training data

(θ) parameterized $R_{\text{emp}}: \mathbb{R}^d \rightarrow \mathbb{R}$ summing up all point losses

Average loss $\bar{R}_{\text{emp}}(f) := \frac{1}{n} \dots$ 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

Usually we find $\hat{\theta}$, the optimal parameter vector, by numerical optimization methods

Stochastic gradient descent

[Geometrically] Find minimum in multidimensional risk surface / loss landscape

But global minimum does not always exist!

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

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

→ local improvement impossible inside very small ε -ball

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

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

- positive eigenvalues
- quadratic form $x^T H x > 0$

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

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

Heads up until we reach local minimum

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

Controlling 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 but 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

→ additional computing memory required

There are improvements for GD such as momentum and ADAM. Implementing those may enable us to smoothenly 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