

Week: 1+2 | Reg+Opt

Week 1
Regression

Linear Regression

Cost function

Week 2
Optimization:

Gradient descent (GD)

Stochastic Gradient descent (SGD)

Mini-Batch Stochastic Gr. Des.
(GPU parallelizable...)

Subgradients:
(non-smooth opt.)
(if \mathcal{L} is not diff. and not necessarily convex)

Either find a relationship between X and Y (Interpretation)
or try to train a model to predict Y (prediction)

$$f_W(x) = x^T W \quad (\text{Linear})$$

$$MSE: \mathcal{L}(w) = \frac{1}{2N} \sum_{n=1}^N (y_n - x_n^T w)^2$$

$$\mathcal{L}(w) = \begin{cases} MAE: \mathcal{L}(w) = \frac{1}{N} \sum_{n=1}^N |y_n - x_n^T w| \\ (= \frac{1}{N} \sum_{n=1}^N (y_n - x_n^T w)) \end{cases}$$

find best $w: w^* \Rightarrow$ Grid Search (not guaranteed, granularity problem)
 \Rightarrow Gradient

$$\frac{1}{2N} \sum_{n=1}^N \|y - X^T w\|^2 = \frac{1}{2N} e^T e \rightarrow \nabla \mathcal{L}(w) = -\frac{1}{N} X^T e$$

$$w^{t+1} = w^t - \delta \nabla \mathcal{L}(w^t) \quad (\delta: \text{Learning rate})$$

$$w^{(t+1)} = w^t - \delta \nabla \mathcal{L}_n(w^{(t)}) \quad (\text{sample one point } n \text{ at random})$$

$$g = \frac{1}{|B|} \sum_{n \in B} \nabla \mathcal{L}_n(w) \rightarrow w^{(t+1)} = w^t - \delta g$$

if $\mathcal{L}(w) = h(q(w))$ i.e. $\mathcal{L}(w) = \frac{1}{N} \sum_{n=1}^N |y_n - x_n^T w|$
where h is not differentiable and q is

Then $g \in \partial h(q(w))$. $\nabla q(w)$ is a subgradient for $\mathcal{L}(w)$ at w

and \rightarrow GD SGD with $g: w^{(t+1)} = w^t - \delta g$

Costs (Lin. Reg: $y = x^T w$)

GD: $O(N \cdot D)$

SGD: $O(D)$

1 step of (S)GD (1-batch)
 $B = N$ (GD)
 $O(|B| \cdot D)$ $\leftarrow B=1$ (SGD)
 $\setminus B$ (H-B SGD)

$$\mathcal{L}(w) = \sum_n \mathcal{L}_n(w)$$

δ : Learning rate
 $\delta = \frac{\epsilon}{t}$ (Cost)
Stopping criteria:
 $\| \nabla \mathcal{L}(w) \| \ll \epsilon$ null K
Convergence:
is only guaranteed when
 $\forall \epsilon > 0 \exists \delta_{\min} = \delta_{\min}$ fixed
depends on the problem