

CS 725

Recap: Algorithm for Gradient Descent

INPUTS: Initial parameter vector w_0 , learning rate $\eta > 0$, threshold $\epsilon > 0$
or step size
maximum # of iterations N_{\max}

OUTPUT: final parameter vector w^*

$w \leftarrow w_0$
for $t = 1 \dots N_{\max}$
 Compute $\nabla_w L(w, \mathcal{D})$
 $w \leftarrow w - \eta \nabla_w L(w)$
 if $(\|\nabla_w L\|_2 \leq \epsilon)$, then break;
return w

$$\begin{aligned}\nabla_w L(w, \mathcal{D}) &= \nabla_w \frac{1}{n} \sum_i (y_i - \underbrace{\hat{y}_i}_{w^T x_i})^2 \\ &= \left(\frac{2}{n} \sum_i (\hat{y}_i - y_i) x_i \right)\end{aligned}$$

EXPENSIVE
Computation involves all n points

// Alternate: if $(\|w_{t+1} - w_t\|_2 \leq \epsilon)$

VARIANTS OF GRADIENT DESCENT

- Full GD : Wt update rule is $W \leftarrow W - \eta \nabla_W L(W, \mathcal{D}_{\text{train}})$
- Stochastic GD (SGD) : Wt update rule is $W \leftarrow W - \eta \nabla_W L(W, \mathcal{D}_{\text{random}})$
Where $\mathcal{D}_{\text{random}} = \{(x_i, y_i)\}$ Where (x_i, y_i) is picked at random from $\mathcal{D}_{\text{train}}$
- (MINI) BATCH GD : Wt update rule is $W \leftarrow W - \eta \nabla_W L(W, \mathcal{D}_{\text{batch}})$
Where $\mathcal{D}_{\text{batch}} = \{(x_i, y_i)\}_{i=1}^B$, Where B is batch size
 $B = |\mathcal{D}_{\text{batch}}|$

Definitions

TRAINING SET:

$\mathcal{D} = \left\{ (x_i, y_i) \right\}_{i=1}^N$, where x_i, y_i are
i.i.d. (independent & identically distributed) points

DEVELOPMENT OR
VALIDATION SET
(HELD OUT SET)

\mathcal{D}_{val} is used to tune hyperparameters

EVALUATION OR
TEST SET

\mathcal{D}_{test} ; not to be used for hyperparam tuning

WEIGHTS or PARAMETERS: w

HYPERPARAMETERS: Predefined values for a chosen model. E.g. $\eta, B, N_{max}, \epsilon, \sigma$,
initial wt vector w_0

TRAIN/DEV/TEST ERROR: $\frac{1}{n} \sum_i (\hat{y}_i - y_i)^2$

PROBABILISTIC MODEL OF LINEAR REGRESSION

Consider $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$. Let each y_i be a noisy target value

defined as: $y_i = f(x_i) + \epsilon_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i$

Where ϵ_i is the noise in the target values that lead to uncertainty in the model.

Let ϵ_i be drawn from a Gaussian distribution with mean 0 and variance σ^2

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$y_i = w^T x_i + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \left[\begin{array}{l} \text{Cov}(\epsilon_i, \epsilon_j) = 0 \\ \text{if } i \neq j \end{array} \right]$$

$$y_i \sim \mathcal{N}(w^T x_i, \sigma^2)$$

$$P(y_i | x_i, w) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left\{ -\frac{(y_i - w^T x_i)^2}{2\sigma^2} \right\}$$

$$P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, \omega) = \prod_{i=1}^n P(y_i | x_i, \omega)$$

(CONDITIONAL) LIKELIHOOD OF DATA

find ω that maximizes the likelihood of observed data \mathcal{D} :

This is MAXIMUM LIKELIHOOD ESTIMATION (MLE)

[frequentist view of ML]

$$\begin{aligned} W_{MLE} &= \operatorname{argmax}_w \prod_i P(y_i/x_i, w) \\ &= \underbrace{\operatorname{argmax}_w \sum_i \log P(y_i/x_i, w)}_{\text{LOG LIKELIHOOD FUNCTION}} \end{aligned}$$

MOTIVATE MLE w/ A COIN TOSSING EXAMPLE

Say we want to estimate the prob of a coin landing on heads, denoted by θ . Consider a dataset of N coin tosses, with N^H heads and N^T tails. What is the max. likelihood estimate of θ ?

Soln:

$$\begin{aligned}\theta_{MLE} &= \arg\max_{\theta} P(\mathcal{D}|\theta) = \arg\max_{\theta} \theta^{N^H} (1-\theta)^{N^T} \\ &= \arg\max_{\theta} \underbrace{N^H \log \theta + N^T \log(1-\theta)}_{LL(\theta)}\end{aligned}$$

$$\frac{\partial L(\theta)}{\partial \theta} = 0 \Rightarrow \frac{N_H}{\theta} - \frac{N_T}{(1-\theta)} = 0$$

$$\Rightarrow \theta_{MLE} = \frac{N_H}{N_H + N_T} = \frac{N_H}{N}$$

Q. Say $\theta \in \{0.4, 0.6\}$

Data: 3 coin tosses, 2H, 1T

What is θ_{MLE} ? 0.6

What is MLE for linear regression?

$$w_{MLE} = \operatorname{argmax}_w \sum_i \log P(y_i | x_i, w)$$

$$= \operatorname{argmax}_w \sum_i \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y_i - w^T x_i)^2}{2\sigma^2} \right\} \right]$$

$$= \operatorname{argmax}_w K - \frac{1}{2\sigma^2} \sum_i (y_i - w^T x_i)^2$$

$$= \operatorname{argmin}_w (y_i - w^T x_i)^2 \quad \text{LEAST SQUARES ERROR}$$