

# CMPUT 267 Basics of Machine Learning

## Winter 2024



# Outline

1. Recap
2. Gradient Descent

# Recap: Estimation

- ▷ True data distribution is  $p_{\text{true}}$ .
  - ▷ Observe dataset  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$  where  $\mathbf{X}_i$  has distribution  $p_{\text{true}}$ .
  - ▷ Goal: Estimate properties of  $p_{\text{true}}$ .
1. Pick a likelihood  $p(\mathcal{D} | \mathbf{w})$  (data generation model).
    - ▷ Pick a distribution class.
  2. Pick a prior that is a conjugate prior  $p(\mathbf{w})$ .
  3. Define objective to get  $\mathbf{w}$ 
    - ▷ MLE:  $\mathbf{c}(\mathbf{w}) = \ln p(\mathcal{D} | \mathbf{w})$
    - ▷ MAP:  $\mathbf{c}(\mathbf{w}) = \ln p(\mathbf{w} | \mathcal{D})$
    - ▷ Bayesian:  $p(\mathbf{w} | \mathcal{D})$

# Recap: Estimation, Conditional models

- ▷ Observe dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  of observations  $\mathbf{x}_i$  and their corresponding targets  $\mathbf{y}_i$ .
- ▷ Goal: Estimate properties of  $P_{\text{true}}(Y|X = \mathbf{x})$ .
- ▷ Our models can be from parametrized families of **conditional distributions**

$$\mathcal{F} = \{f(\mathbf{y} | \mathbf{x}; \theta) \mid \theta \in \mathbb{R}^k\}$$

1. Pick a likelihood  $p(\mathbf{y} | \mathbf{x}, \mathbf{w})$  (data generation model).
  - ▷ Pick a distribution class, prior.
  - ▷ MLE:  $\mathbf{c}(\mathbf{w}) = \sum_{i=1}^n \ln p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w})$

Example:

$$p(y|x) = N(xw_1, \exp(xw_2))$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \in \mathbb{R}^2$$

MLE

$$\min_{\vec{w} \in \mathbb{R}^2} \sum_{i=1}^n -\ln p(y_i | x_i; \vec{w}) \leftarrow \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - y_i)^2}{2\sigma^2}\right)$$

$$-\ln p(y_i | x_i; \vec{w}) = -\ln \frac{1}{\sqrt{2\pi \exp(x_i w_2)}} - \ln \left( \exp\left(-\frac{(x_i w_1 - y_i)^2}{2 \exp(x_i w_2)}\right) \right)$$

$$= \ln(\sqrt{2\pi}) + \frac{1}{2} x_i w_2 + \frac{(x_i w_1 - y_i)^2}{2 \exp(x_i w_2)}$$

$$c(w) = \sum_{i=1}^n c_i(w) \quad \text{where} \quad c_i(w) = \frac{1}{2} \left[ x_i w_2 + \frac{(x_i w_1 - y_i)^2}{\exp(x_i w_2)} \right]$$

$$\frac{\partial c_i(w)}{\partial w_1} = \frac{(x_i w_1 - y_i) x_i}{\exp(x_i w_2)}$$

$$\frac{\partial c_i(w)}{\partial w_2} = \frac{x_i}{2} + \frac{(x_i w_1 - y_i)^2}{2} \frac{\partial \exp(-x_i w_2)}{\partial w_2}$$

$$\begin{bmatrix} \frac{\partial c(w)}{\partial w_1} \\ \frac{\partial c(w)}{\partial w_2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{\partial c(w)}{\partial w_1} = \sum_{i=1}^n \frac{x_i (x_i w_1 - y_i)}{\exp(x_i w_2)} = 0$$

$$\frac{\partial c(w)}{\partial w_2} = \frac{1}{2} \sum_{i=1}^n \left[ x_i - x_i (x_i w_1 - y_i) \exp(-x_i w_2) \right] = 0$$

→ No closed form; we need to use gradient descent

GD

Initialize  $\vec{w}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$t=0 \quad \vec{g}_0 = \nabla c(\vec{w}_0) = \begin{bmatrix} \frac{\partial c(\vec{w}_0)}{\partial w_1} \\ \frac{\partial c(\vec{w}_0)}{\partial w_2} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n -x_i y_i \\ \sum_{i=1}^n x_i + x_i y_i \end{bmatrix} \leftarrow \vec{g}_{0,1}$$

$$w_{t+1} \leftarrow w_t - \eta \vec{g}_t$$

stopping condition: until  $\|\vec{g}_t\|_2^2 = \sum_{j=1}^2 g_{t,j}^2 \leq \epsilon$  (threshold)

Problem: GD can be very expensive for large dataset

Stochastic Gradient Descent : take a stochastic approximation to  $\nabla C(w)$

$$\nabla C(w) = \frac{1}{n} \sum_{i=1}^n \nabla C_i(w)$$

→ take a subsample of population

sample average

$\nabla C_k(w)$  :  $k$ : random sample

$$= \frac{1}{|k|} \sum_{i \in k} \nabla C_i(w)$$

unbiased estimate of  $\nabla C(w)$

sample size of  $k$

$$E[\nabla C_k(w)] = \sum_{i=1}^n \underbrace{p(k)}_{1/n} \nabla C_i(w) = \frac{1}{n} \sum_{i=1}^n \nabla C_i(w) = \nabla C(w)$$

unbiased

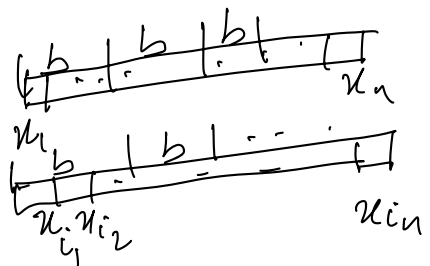
But variance may be high.

We might be able to reduce the variance by selecting a batch of samples instead of a single one at each iteration.

SGD with minibatches

batch size of  $b \geq 1$

shuffle →



$$\nabla C(w) \approx \frac{1}{b} \sum_{j=1}^b \nabla C_{ij}(w_t)$$

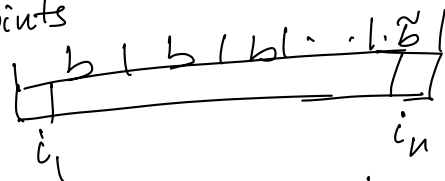
$g_t$

$$w_{t+1} = w_t - \eta_t g_t$$

1. Initialize  $\vec{w}_0$

For  $t=1$  to  $T$  ( $T$ : nb of epochs)

shuffle the order of datapoints



for  $k=0, \dots, \lfloor \frac{n}{b} \rfloor$

select  $k^{\text{th}}$  minibatch

$\lfloor \frac{n}{b} \rfloor$  batches  
of size  $b$

compute  $\hat{g}_t$  from this minibatch

$$\hat{g}_t = \frac{1}{b} \sum_{i=0}^{b-1} \nabla c_{i_n}(w_t)$$

update the parameter guess

$$\vec{w}_{t+1} = \vec{w}_t - \eta_t \hat{g}_t$$

$t = t+1$

Computational cost

$w \in \mathbb{R}^d$

$$\text{GD: } w_{t+1} = w_t - \eta_t \frac{1}{n} \sum_{i=1}^n \nabla c_i(w_t)$$

$$\text{SGD } w_{t+1} = w_t - \eta_t \frac{1}{b} \sum_{k=1}^b \nabla c_{i_k}(w_t)$$

cost of one update  $O(nd)$

$O(bd)$   
 $b \ll n$

Full cost

$$O(T_{\text{GD}} n \cdot d)$$

$T_{\text{GD}}$ : #updates  
for GD

$$O(T_{\text{SGD}} \cdot bd)$$

$T_{\text{SGD}}$ : #updates for SGD

$$T_{\text{SGD}} \cdot bd < T_{\text{GD}} \cdot nd \rightarrow T_{\text{SGD}} < T_{\text{GD}} \cdot \frac{n}{b}$$

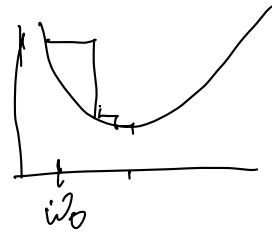
Ex.  $n: 1M$     $b: 32$     $T_{\text{SGD}} < 30,000 T_{\text{GD}}$

# Selection of step size for SGD

- Fixed step size
  - crude, may take too many iterations to get to optimal point.

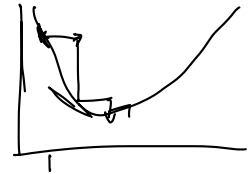
- Heuristics

$$\eta \approx \frac{1}{t}$$



- Gradient dependent

$$\eta_t = \frac{1}{1 + |g_t|}$$



$$w_{t+1} = w_t - \frac{g_t}{1 + |g_t|}$$

$$C(w) = \frac{1}{2} w^2 \quad \nabla C(w) = w$$

If  $w_0$ :

$$w_0 = 1.0 \quad \eta \approx \frac{1}{t}$$

$$\eta = \frac{1/2}{t+1}$$

$$w_1 = w_0 - \eta_0 w_0$$

$$= 1 - \frac{1}{2}(1)$$

$$= 1/2$$

$$\eta_t = \frac{1/2}{2} = 1/4 \quad w_2 = w_1 - \eta_1 w_1$$

$$= \frac{1}{2} - \frac{1}{4} \left( \frac{1}{2} \right)$$

$$= 0.375$$

$$\eta_2 = 0.16$$

$$w_3 = 0.3125$$

$$\eta_3 = 0.125$$

$$w_4 = 0.273$$

$$\eta_4 = 0.1$$

$$w_5 = 0.246$$

Gradient-dep.

$$\eta_t = \frac{1}{1 + |g_t|}$$

$$w_0 = 1 \quad \eta_0 = \frac{1}{2}$$

$$w_1 = w_0 - \eta_0 w_0$$

$$= 1/2$$

$$\eta_1 = \frac{1}{1 + |g_1|} = \frac{2}{3}$$

$$w_2 = \frac{1}{2} - \frac{2}{3} \left( \frac{1}{2} \right)$$

$$0.16$$

$$\eta_2 = 0.85$$

$$w_3 = 0.02$$

$$\eta_3 = 0.98$$

$$w_4 \approx 5 \times 10^{-4}$$

$$\eta_4 = 0.99$$

$$w_5 \approx 3 \times 10^{-7}$$



SGD: noisy gradient  $\tilde{\nabla} c(w_t) = w_t + \epsilon_t$   
 $\epsilon_t \sim \mathcal{N}(0, 0.2)$

$$\eta_t = \frac{1/2}{t+1}$$

- $w_0 = 1.0$
- $w_1 = 0.58$
- $w_2 = 0.44$
- $w_3 = 0.35$
- $w_4 = 0.009$
- $\vdots$

$$w_{1000} \approx 0.003$$

$$E\{\tilde{\nabla} c(w_t)\} = \nabla c(w_t)$$

$$\eta_t = 1/(t+1) |g_t|$$

- $w_0 = 1.0$
- $w_1 = 0.47$
- $w_2 = 0.09$
- $w_3 = -0.06$
- $w_4 = -0.24$
- $w_5 = 0.03$
- $\vdots$

$$w_{10000} \approx 0.05$$

Adagrad

$$\eta_t = \frac{1}{1 + \bar{g}_t}$$

$$\bar{g}_t = \bar{g}_{t-1} + |g_t|$$

$$\bar{g}_0 = |g_0|$$

$$\underline{\bar{g}_{t+1} \geq \bar{g}_t}$$

$$\underline{\eta_{t+1} \leq \eta_t}$$

slowly decreasing