

Q3

1): $E_I(\frac{1}{m} \sum_{i \in I} a_i)$

$$= \frac{1}{m} E(a_{i1} + a_{i2} + \dots + a_{im})$$

$$= \frac{1}{m} [E(a_{i1}) + \dots + E(a_{im})]$$

$$= \frac{1}{m} \left[\sum_{i=1}^n a_i \left(\frac{1}{n}\right) + \dots + \sum_{i=1}^n a_i \left(\frac{1}{n}\right) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n a_i$$

2): $E_I(\nabla L_I(x, y, \theta))$

$$= E_I \left[\frac{1}{m} \nabla \sum_{i \in I} L(x, y, \theta) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \nabla L(x^{(i)}, y^{(i)}, \theta) \quad (\text{by (1)})$$

$$= \nabla \left[\frac{1}{n} \sum_{i=1}^n L(x^{(i)}, y^{(i)}, \theta) \right]$$

$$= \nabla L(x, y, \theta)$$

3): Mini-batch method produces an unbiased estimator of true gradient

4): a): $L(x, y, w) = (y - w^T x)^2$

Assume $w: d \times 1$; $x: d \times N$; $y: 1 \times N$

$$\nabla_w L(x, y, w)$$

$$= \nabla_w (y y^T + w^T x x^T w - 2 w^T x y^T)$$

$$= 2 x x^T w - 2 x y^T$$

5): cos-similarity is more meaningful

Since gradient computed has large values. Even with cos-similarity over 99%, square distance can still be quite large

6) ρ/bb

