

Answers sheet - Final Exam DSE230, June 2016

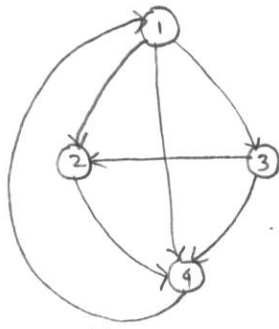
Name: RAJESH BONDUQUA

PID: A53108479

Q1 (3pt): ①. 2. 3. ④.	Q2 (3pt): 1. 2. ③. 4.
Q3 (10pt): Method 1 is faster. Method 1 has only one reduce and hence can be completed in 1 Pass. Method 2 has 2 reduce calls, one for mean calculation another for variance calculation from mean. This does 2 scans on data and hence is slower than Method 1.	
Q4 (6pt): <code>S = genderRDD.sample(False, 0.0001) # Sample RDD to get 1000 elements</code> <code>freq = S.countByValue() # Returns dictionary of frequencies</code> <code>male_pct = freq['m'] * 100.0 / (freq['m'] + freq['f']) # Male pct</code> <code>female_pct = freq['f'] * 100.0 / (freq['m'] + freq['f']) # Female pct</code>	
Q5 (14pt): <code>grp = graphRDD.groupByKey().mapValues(list) # Grp edges based on starting node</code> <code>t = grp.top(1, Key=lambda x: len(x[1])) # Get the top node and its connectives</code> Print "top node = ", <code>t[0]</code> , "out-degree = ", <code>len(t[1])</code>	
Q6 (5pt): 1. ②. 3. 4.	
Q7 (6pt): $+0.5 - 0.7 - 0.2 - 0.1 = -0.5$. Sum = -0.5 and hence it is classified as "-1"	
Q8 (15pt): % of variance explained by top 1 eig vector = $\lambda_1 / (\lambda_1 + \lambda_2 + \lambda_3)$ % of variance explained by top 2 eig vector = $(\lambda_1 + \lambda_2) / (\lambda_1 + \lambda_2 + \lambda_3)$ % of variance explained by top 3 eig vector = $(\lambda_1 + \lambda_2 + \lambda_3) / (\lambda_1 + \lambda_2 + \lambda_3) = 1.0$ approximation of $x = \mu + ((x - \mu) \cdot v_1) v_1 + ((x - \mu) \cdot v_2) v_2$ approximation using $v_3 = x - ((x - \mu) \cdot v_3) v_3$	
Q9 (8pt): $\frac{\delta}{\delta x} x - k = 1$ if $x > k$ else -1 if $x < k$ $\frac{\delta}{\delta x_1} F(\vec{x}) = 1$, $\frac{\delta}{\delta x_2} F(\vec{x}) = -1$, $\frac{\delta}{\delta x_3} F(\vec{x}) = -1$, $\frac{\delta}{\delta x_4} F(\vec{x}) = 1$ \therefore gradient = $(1, -1, -1, 1)$	
Q10 (10pt): 1. ②. ③. 4. ⑤.	

$$x = ((x - \mu) \cdot v_1) v_1 + ((x - \mu) \cdot v_2) v_2 + ((x - \mu) \cdot v_3) v_3 + \mu$$

Q5



For this example RDD has $[(1,2), (1,4), (1,3), (2,4), (3,2), (3,4), (4,1)]$ (say this is graph RDD)

`graphRDD.groupByKey().mapValues(list)` results $[(1, [2, 4, 3]), (2, [4]), (3, [2, 4]), (4, [1])]$

`rx.top(1, Key = lambda x: len(x[1]))` will return $(1, [2, 4, 3])$

$x[0]$ is top node

$\text{len}(x[1])$ is out degree