10/25

$$\begin{bmatrix} a_1, a_2 \\ a_3, a_4 \end{bmatrix} 32 \text{ bits} \times 4$$

$$= 128 \text{ bits}$$

$$\begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \cdot \begin{bmatrix} b_1 & b_2 \\ b_3 & b_4 \end{bmatrix} =$$

2×2

$$\begin{bmatrix} a_1 \cdot b_1 + a_2 \cdot b_3, & a_1 \cdot b_2 + a_2 \cdot b_4 \\ a_3 \cdot b_1 + a_4 \cdot b_3, & a_3 \cdot b_2 + a_4 \cdot b_4 \end{bmatrix}$$



→ | 1 | 4 | 7 | 2 | 5 | 8 | 3 | 6 | 9 |

column major

↳ If you ignore the sparsity
you may store the matrix
in either row-major order
or column-major order.

CSR = compressed sparse rows.
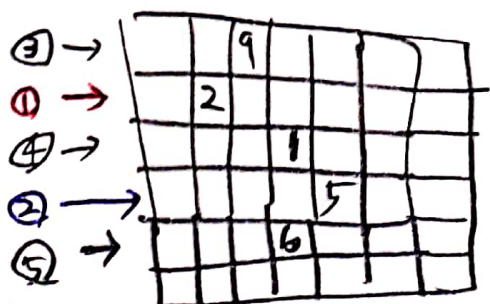
↳ store the column of each non-zero element, and which column indices belong to each row.



| 0 | 2 | 3 | 3 | 3 | 6 | 6 | 7 |

0+0 = 0 ... =2 0+2

2+1 = 3

0
1
2
3

Val = | 3 | 7 | 8 | 2 | 5 | 9 |

Col = | 1 | 2 | 2 | 0 | 1 | 3 |

Row = | 0 | 0 | 0 | 6 | ... |

Val = | 3 | 7 | 8 | 2 | 5 | 9 |
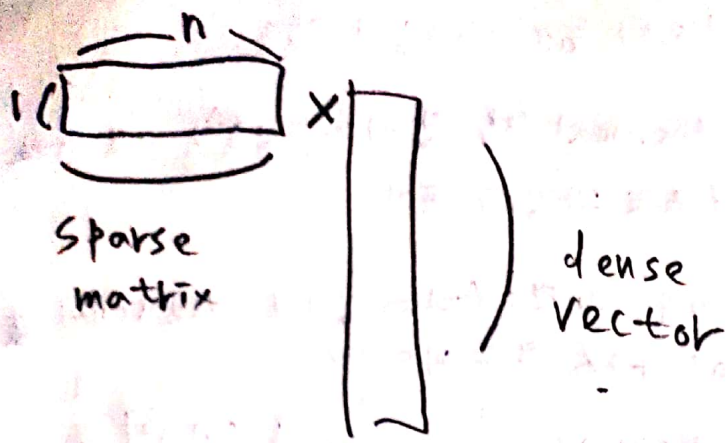0   2   3   5   6

Popular format = ⌐ COO, CSR ⌐

Row = | 1 | 3 | 0 | 2 | 4 |
      ①   ②   ③   ④   ⑤

data = | 2 | 5 | 9 | 1 | 6 |
       ①   ②   ③   ④   ⑤



③ →
① →
④ →
② →
⑤ →

<1>

Sparse matrix × dense vector

$k$ × (M×K) → M, K ... (k×1) ... B matrix (dense)

Accumulator

⇒ vector-wise operation.

grouping multiple vectors lead to more number of wasted computation.



2 : grouping multiple Vectors.

● 위치다르네가 dense vector가 이전보다 계산 낭비↑↑.

---

meta data



R

$4 = c/2$ 「2-bits indices」

<2>

---

<Pruning Initialization>

traing → [Pruning] → retraining

Folk knowledge [without an initial training]

① The pruned model cannot be optimized well

② final model does not generalize well

---

「Can prune at the initial stage without performance drop.」

---

↳ There are other methods called Grasp/synFlow, but SNIP perform best.
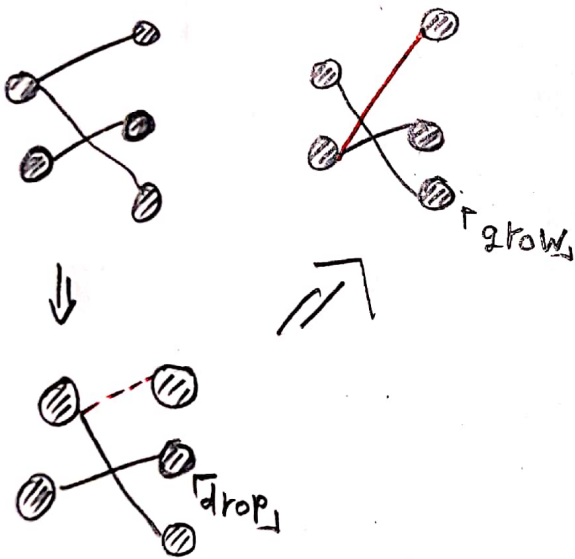
<3>

layerwise sparsity.

taylor approximation

$\hookrightarrow L(\tilde{\theta}) \simeq L(\theta) + (\tilde{\theta}-\theta)^T G_\theta + \frac{1}{2}(\tilde{\theta}-\theta)^T H_\theta \cdot (\tilde{\theta}-\theta)$

sparse training = Limitation
: Peak memory is still the same as dense.

⟱

↗

「drop」

「grow」

Sparsity ↓↓ ⟹ test accuracy.
↑↑

「good layerwise sparsity & schedule」

pruing

① minimizing the loss after pruning

② maximizing the re-train ability

① = 더 중요 (예전)

↓

(2) is viewed as the most important decision criterion.