

Object localization.

Softmax를 통해 나온 레이블 $y =$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

(p_c 는 물체가 있는지에 대한 확률)
($p_c=0$ 이면 물체가 없는 배경임을 의미.)

물체(분류)

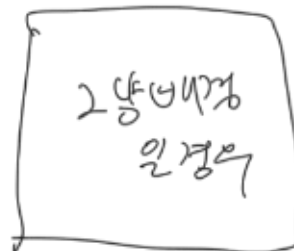
대한 클래스값

b_x, b_y, b_w, b_h 는 경계상자에 대한 변수값이다.



자율주행차라면, 예) $\begin{cases} c_1 = 보행자 \\ c_2 = 차 \\ c_3 = 오토바이 \end{cases}$

\therefore 이 과정은 이미지가 단 한개의 물체만 포함한다는 가정



분류레이블 y

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

p_c 가 0이므로 경계상자
 b_x, b_y, b_h 와 물체클래스
분류를 할 필요가 없어진다.

'don't care'

이렇게 소프트맥스층을 통해 얻은 분류레이블 y 를 가지고

신상차-2 계신상차-1000

이제부터 MV 회귀 분석

Mean Squared Error로 보면

$$L(\hat{y}_i, y_i)$$

실제값 예측값

$$E = \sum (\hat{y}_i - y_i)^2 \rightarrow \begin{array}{l} y \text{의 평균} \\ \text{실제 } y = 1 \text{ 이고} \\ P_c \text{ 역시 } 1 \text{ 이고 } 0 \\ \vdots \end{array}$$

☞ P_c 즉, 물체가 있는지 인식한 확률에 따라 손실함수 값이 다르다.

i) $P_c = 0$ (물체가 없는 경우)

ii) $P_c = 1$ (물체가 있음)

그러면 $L(\hat{y}_i, y_i)$ 는

$\therefore (\hat{y}_i - y_i)^2$ 이 된다.

그러면 물체가 있으면

두(에) 출력값 신경망의

$$\sum_i (\hat{y}_i - y_i)^2 \text{ 로는 물체가 없는 경우에 대한}$$

오류 제곱의 합

이렇게 나누어서 손실함수 계산.

실제로는 자세히 보면

$$\begin{bmatrix} P_c \rightarrow \text{로지스틱 확률 손실함수 사용} \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \text{ 은 MSE 사용}$$

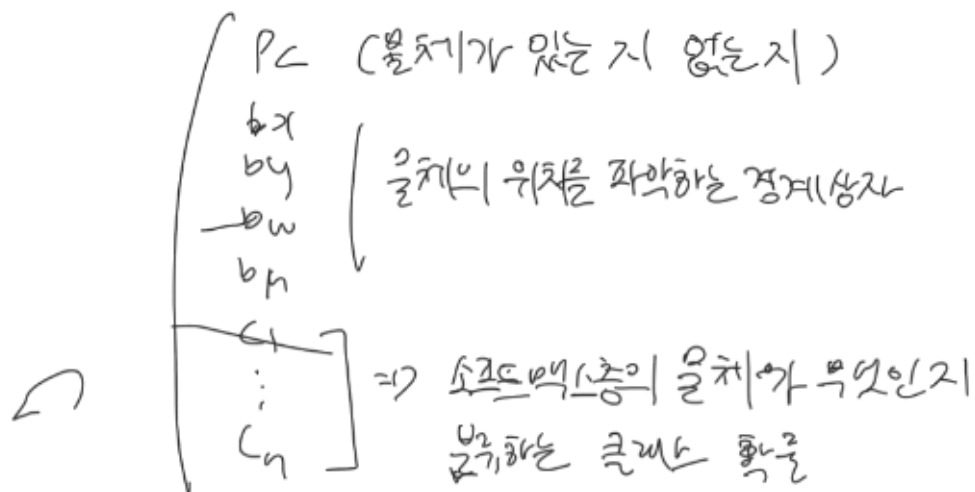
$$\begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} \text{ 은 로지스틱 손실함수 사용}$$

$\therefore P_c$ 즉 y_i 의 값을 얼마나 정확히 계산하는가에 따라 손실함수 값이 좌우된다.

\Rightarrow 물체로 분류 리제이션 \approx 물체를 인식함수를 넣어

물체의 위치가 어디 있는지 파악해볼
(경계 상자)을 이용한다.

이름은 소프트맥스층을 이용한 분류레이블을 출력하려고



이를 출력하고 이 값 중 $P_C = 0$ 이 나오면

손실함수를 계산 $\rightarrow P_C$ 가 0이면 나머지
 7개 항목
 don't care 한데,
 (물체가 없으니까)

* 특징점 탐지
 = landmark detection

\rightarrow ex) 얼굴 사진에서 눈을 탐지하고자 한다면



\downarrow
 CONV Network

\downarrow
 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$ softmax

각각 눈꼬리 양 끝점 중 4개의
 $\begin{cases} l_{1x}, l_{1y} \\ l_{2x}, l_{2y} \\ l_{3x}, l_{3y} \\ l_{4x}, l_{4y} \end{cases}$
 을 출력 레이블로
 나타내면 된다.

$\downarrow \downarrow \downarrow \dots \downarrow$
 $x_1, x_2, x_3, \dots, x_n$

* 이것이 실전에서 감정을
 인식하는 것의 기본적인 설계단위이다.

⇒ 특징점 탐색 ⇒ 특징을 통해 인식하고자 하는 특징점들의
 각 좌표를 출력한다.

(특징점은 다른 이미지에 대해서도 항상 동일해야 한다.)



즉, 레이블의 순서는 다른 이미지에 대해서도 항상 동일해야 한다.

* 물체 인식 (Object detection) 알고리즘 설계방법.

가장 먼저 label training set을 만든다.

Training set :



이 레이블 표현식으로
 CNN을 훈련시킬 수 있다.

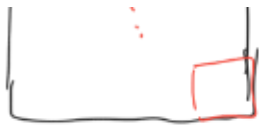
Crop된 이미지를 사용해서
 학습하게 train



* sliding windows detection



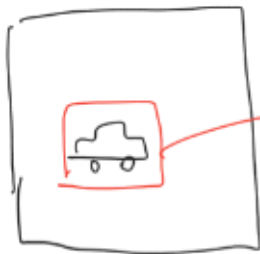
항상 고정된 크기의 이 작은 시각영역을 입력



□: window값은 이미지의 작은영역을 잡아서
합성곱 신경망에 통과시켜 예측값을 구한다.

이미지의 잘려진 부분에 대한 합성곱 신경망을 통과한
예측값은 각 위치에서 0~1의 값으로 분류된다.

* 한번이 sliding window가 모두 끝나면
그다음은 더 큰 □을 CONVNet에 통과시킨다



이렇게 차가 지나면 1을 출력함

이것은 sliding window detection 이라한다.

* sliding window 조출의

단점은 계산 비용이다.

→ 수많은 작은영역 계산에 대한

* 해결: 슬라이딩 윈도우 알고리즘을 합성곱을 통해
구현해낼수 있다.

* Turning FC Layer into Convolutional layers



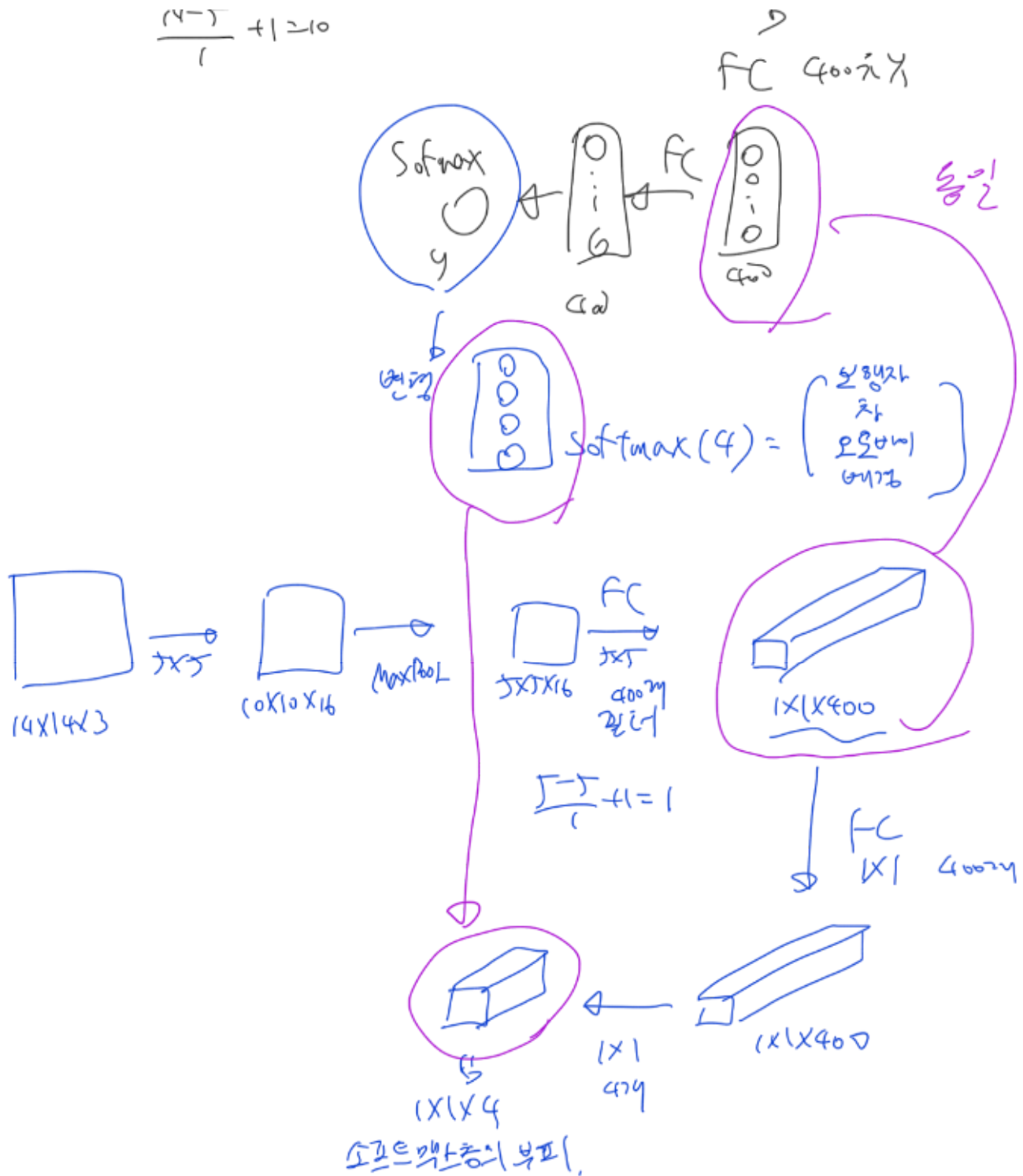
$f: 5 \times 5$
1674



MaxPool
2x2



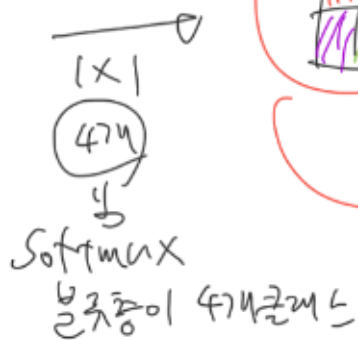
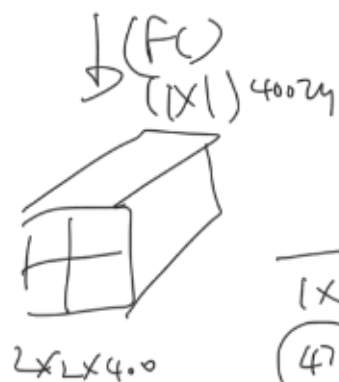
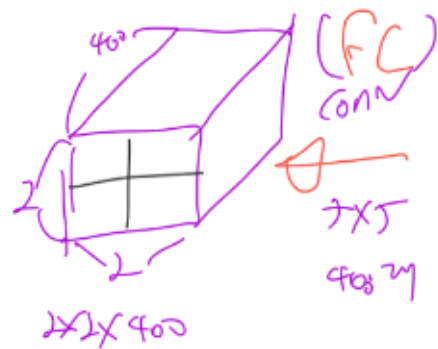
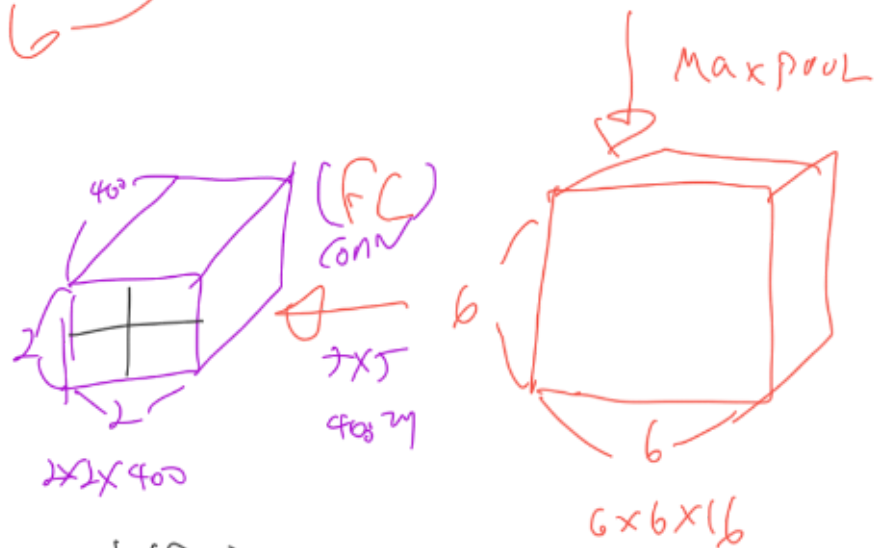
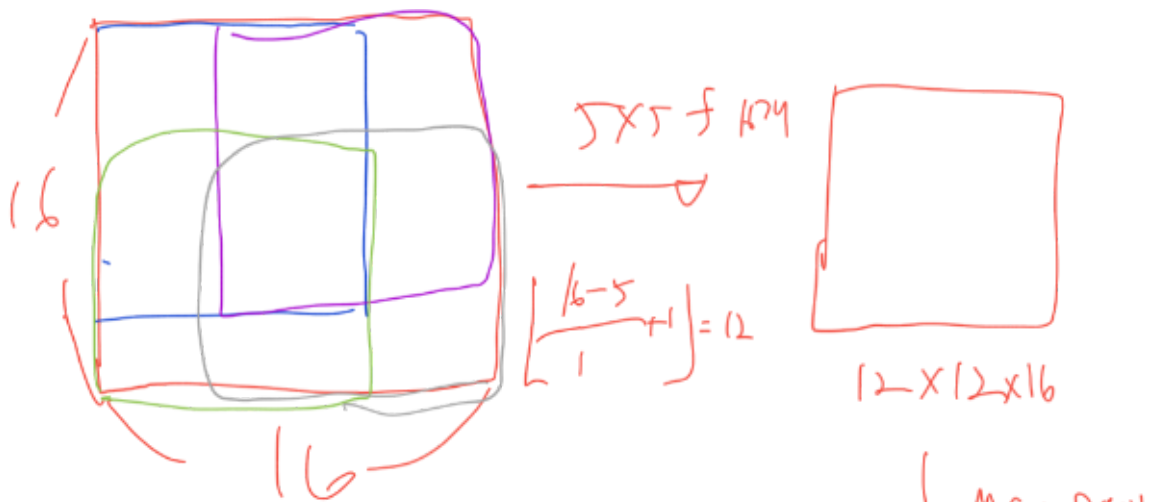
$$\frac{(N-1)}{1} + 1 = 10$$



FC를 이용해 합성곱 층을 구현해낸 것.

뉴네이빙 인도를 합성곱 Net을 통해 구현함의 장점은 계산 변수를 공유한다는 것.

ex) 4개의 원소를 뉴네이빙하여 물체 검출에서



이 4개의 영역이

클래싱 원도 4번하여 나온 값이 24.
즉 4번 계산한것을 합성곱을 통해
1번 계산으로 다 알 수 있죠.

마지막 수정: 오후 5:08