

THIAGO MALTA COUTINHO - 2014123335

- ① Porque segmentação semântica atribui uma classe a cada pixel e detecção de objetos atribui uma classe a um conjunto de pixels definidos pelo bounding box.
- ② Classificação de imagem tem o objetivo de dizer se há ou não uma classe na imagem. Detecção de objetos determina a localização do objeto na imagem e o classifica. Segmentação semântica classifica cada pixel da imagem. Segmentação de instância determina a localização do objeto na imagem e classifica os pixels do objeto.
OBS: Detecção de objetos determina a localização utilizando uma bounding box, nem todos os pixels da caixa são pertencentes ao objeto classificado.
- ③ Montar a rede e utilizar uma camada softmax na saída. O intervalo $[0, 1]$ será mapeado para um valor na escala de cinza $[0, 255]$. Cada pixel recebe uma probabilidade e um valor de cinza no mapa final.
- ④ Efetuar recorte na imagem original e passar para a rede. O pixel central é o pixel que será classificado. Essa abordagem é custosa e alguns pixels de borda serão perdidos caso não se use padding.
- ⑤ Uma rede Fully connected com um neurônio para cada pixel não mantém corretamente a informação espacial dos pixels, deteriorando o desempenho da rede.
- ⑥ Downsampling é feito ao aplicar filtros convolucionais. O objetivo é extrair/gerar características complexas que representam bem o problema e o tornar separável/"classificável". Upsampling é feito com o objetivo de levar as features complexas do downsampling e a percoar com o tamanho original da imagem, tornando a

informação utilizada.

7

$$\begin{bmatrix} 5 & 6 \\ 3 & 2 \end{bmatrix}$$

SAÍDA: 4x4

a)
$$\begin{bmatrix} 5 & 5 & 6 & 6 \\ 5 & 5 & 6 & 6 \\ 3 & 3 & 2 & 2 \\ 3 & 3 & 2 & 2 \end{bmatrix}$$

b)
$$\begin{bmatrix} 5 & 0 & 6 & 0 \\ 0 & 0 & 0 & 0 \\ 3 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

c)
$$\begin{bmatrix} 3 & 1 & 1 & 4 \\ 2 & 1 & 3 & 2 \\ 3 & 7 & 1 & 7 \\ 3 & 2 & 2 & 4 \end{bmatrix}$$

MAX-POOL

$$\begin{bmatrix} 3 & 4 \\ 7 & 7 \end{bmatrix}$$

$$\begin{bmatrix} 5 & 6 \\ 3 & 2 \end{bmatrix}$$

UNPOOLING

$$\begin{bmatrix} 5 & 0 & 0 & 6 \\ 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 2 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

d)
$$\begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} \rightarrow 5 \cdot \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 5 \\ 10 & 15 \end{bmatrix}$$

$$\rightarrow 6 \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 6 \\ 12 & 18 \end{bmatrix}$$

$$\rightarrow 3 \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 3 \\ 6 & 9 \end{bmatrix}$$

$$\rightarrow 2 \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 2 \\ 4 & 6 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 5+0 & 6 \\ 10+0 & 15+12+3+0 & 18+2 \\ 6 & 9+4 & 6 \end{bmatrix}$$

||

$$\begin{bmatrix} 0 & 5 & 6 \\ 10 & 30 & 20 \\ 6 & 13 & 6 \end{bmatrix}$$

THIAGO MALTA COUTINHO - 2014123335

$$\textcircled{1} \quad y = \begin{bmatrix} 1 \\ 0,4 \\ 0,8 \\ 0,2 \\ 0,2 \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{array}{l} \rightarrow \text{OBJECT?}(P_c) \\ \rightarrow b_x \\ \rightarrow b_y \\ \rightarrow b_h \\ \rightarrow b_w \\ \text{CLASSES ONE-HOT} \end{array}$$

$$\textcircled{2} \quad y = [0 \ ? \ ? \ ? \ ? \ ? \ ? \ ?]^T$$

Os pesos da rede não serão atualizados.

$$\textcircled{3} \quad L(w) = \begin{cases} \sum_i (y_i - \hat{y}_i)^2, & \text{se } y_{\perp} = 1 \\ (y_{\perp} - \hat{y}_{\perp})^2, & \text{se } y_{\perp} = 0 \end{cases}$$

$\textcircled{4}$ Como é uma classificação binária, a última camada será uma regressão logística. Apenas b_x e b_y são necessários, visto que as dimensões da lata são fixas.

$$y = [P_c \ b_x \ b_y]^T$$

$\rightarrow (x, y) \rightarrow$ centro da bounding box/lata.
 \rightarrow Probabilidade de haver uma lata.

$\textcircled{5}$ Como é uma imagem 2D, cada ponto tem coordenada (x, y) , logo 2N.

$\textcircled{6}$ É necessário fornecer vetores ONE-HOT e as bounding boxes com coordenadas (b_x, b_y) associadas ao centro e b_w e b_h associadas à largura da caixa e altura, respectivamente.

⑦ Aumentar o stride reduz o custo computacional pois menos valores são processados mas diminui a precisão, pois menos informação é processada. As alternativas são: convolução dilatada, sliding window convolucional e region proposal.

⑧ R-CNN faz a proposição de regiões e então o feature extraction.

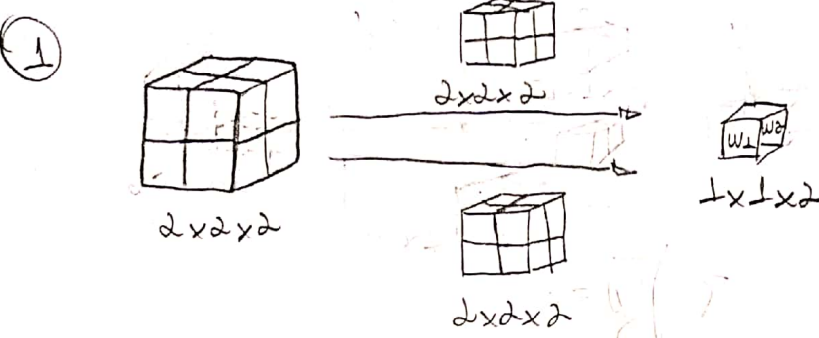
Já o Fast R-CNN faz o feature extraction na imagem completa e region proposal na imagem original. O ROI pool é utilizado para alinhar o mapa com a ROI e transformar as dimensões para o tamanho de saída necessário. A velocidade do forward é aumentada significativamente.

⑨ A faster R-CNN utiliza uma rede CNN para propor as regiões de interesse, diferente da Fast R-CNN, que utiliza Selective Search.

⑩ ROI pooling com alinhamento produz um mapa de ativação de saída compatível com a entrada de um classificador, podendo ser uma CNN. Ele padroniza a dimensão de saída para diferentes tamanhos de ROI.

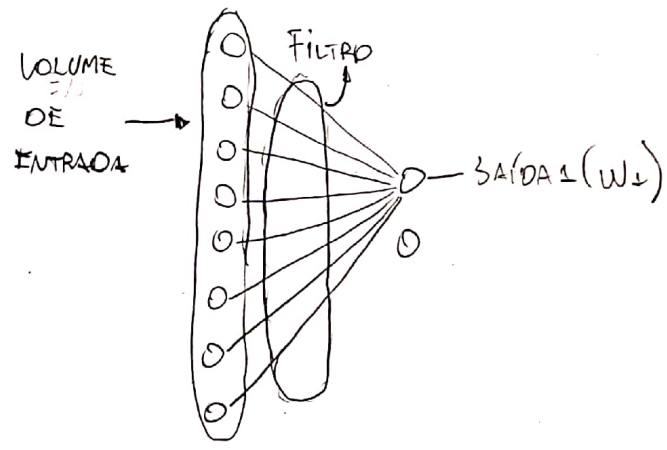
LISTA TEÓRICA - DETECÇÃO / VOLO

THIAGO MARTA LOUTINHO - 2014123335



$$W_k = g\left(\sum_{i,j,c=1}^2 w_{i,j,c} F_{i,j,c,k} + b_k\right)$$

São 8 entradas, $(2 \times 2 \times 2) \times 2$ pesos, 2 saídas. Cada filtro é equivalente a um neurônio e os valores $F_{i,j,c,k}$ são as conexões da entrada com o filtro/neurônio correspondente. É uma camada F_c com 2 neurônios.



② $I_{OU} = \frac{A \cap B}{A \cup B}$? $A \cap B = 1$ $\Rightarrow I_{OU} = 1/9$
 $A \cup B = 9$

③ $P(x) \leq 0,4 \rightarrow$ DESCARTA
 $I_{OU} > 0,5 \rightarrow$ SOBREPosição

5 CAIXAS NO TOTAL

- $P(MOTORCYCLE) = 0,58 \rightarrow OK$
 - $P(PEDESTRIAN) = 0,98 \rightarrow OK$
 - $P(CAR1) = 0,26 \rightarrow X$
 - $P(CAR2) = 0,63 \rightarrow X$
 - $P(CAR2) = 0,73 \rightarrow OK$
 - $P(TREE) = 0,46 \rightarrow OK$
 - $P(TREE) = 0,76 \rightarrow OK$
- $I_{OU} > 0,5 \rightarrow OK$
 $I_{OU} > 0,5 \rightarrow X$

④ $P(x) \notin 0,6$

$P(\text{MOTORCYCLE}) = 0,58 \rightarrow x$

$P(\text{CAR 2}) = 0,62 \rightarrow \cancel{x}$
 $P(\text{CAR 2}) = 0,73 \rightarrow \text{OK}$ } IOU > 0,5 OK

3 CAIXAS NO TOTAL

$P(\text{CAR 1}) = 0,26 \rightarrow x$

$P(\text{PEDESTRIAN}) = 0,98 \rightarrow \text{OK}$

$P(\text{TREE}) = 0,46 \rightarrow x$

$P(\text{TREE}) = 0,74 \rightarrow \text{OK}$

⑤ Objetos distintos sobrepostos

Bounding box com maior probabilidade não é o mais adequado.

⑥ Apenas a célula com anchor que identifica o objeto.

⑦ $\left\{ [P_c \ b_x \ b_y \ b_h \ b_w \ \text{ONE-HOT}_{20}] \times 5 \right\} (\text{grid}) = (25 \times 5) \times (19 \times 19)$

⑧ YOLO nem sempre produz os melhores resultados. Ele possui dificuldades com proporções distorcidas, limitações de tamanho de grid.

⑨ Porque a sobreposição entre os objetos resultará em probabilidade baixa para o objeto sobreposto, que será filtrado pelo non-max suppression.

⑩ Objetos de tamanhos diferentes e em posições diversas serão detectados.

⑪ Manualmente \rightarrow simplicidade na implementação, mas com risco de erro.
 k-means \rightarrow melhor ajuste das anchor boxes.

LISTA TEÓRICA - AUTOENCODERS

THIAGO MALTA COUTINHO - 2014123335

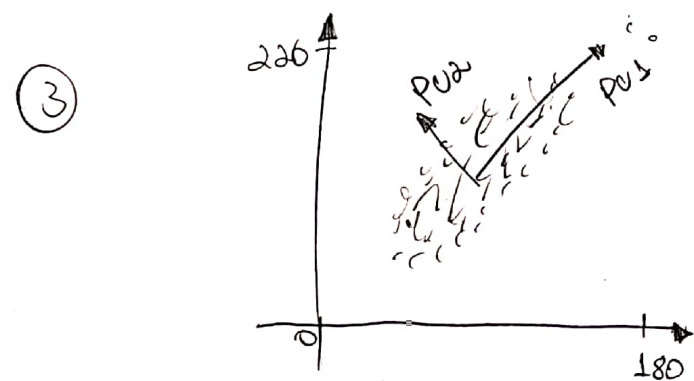
① Aprendizado supervisionado possui dados rotulados, já o não-supervisionado, não.

Supervisionado: classificação, regressão.

Não-supervisionado: clustering, redução de dimensionalidade.

② $\Sigma \vec{v} = \lambda \vec{v} \quad ; \quad \vec{v} = \{1, 1, 1\} / \sqrt{3}$

$$Z = \begin{bmatrix} 1 & p & p \\ p & 1 & p \\ p & p & 1 \end{bmatrix} \quad \Sigma \vec{v} = \begin{bmatrix} 1+p+p \\ 2p+1 \\ 2p+1 \end{bmatrix} \rightarrow (2p+1) \rightarrow \text{AUTOVALOR}$$



④ A ideia principal é aprender codificar a variável de entrada em um espaço latente e decodificá-la com a menor perda de informação possível. Nesse processo o espaço latente é aprendido e pode ser utilizado para diversas aplicações como feature engineering. O autoencoder tem a estrutura de encoder e decoder, ambas são redes neurais.

⑤ Detecção de anomalias baseada no erro de reconstrução.

- Feature engineering a partir do espaço latente aprendido.
- Pré-treinamento para uma rede que será concatenada na saída do decoder.

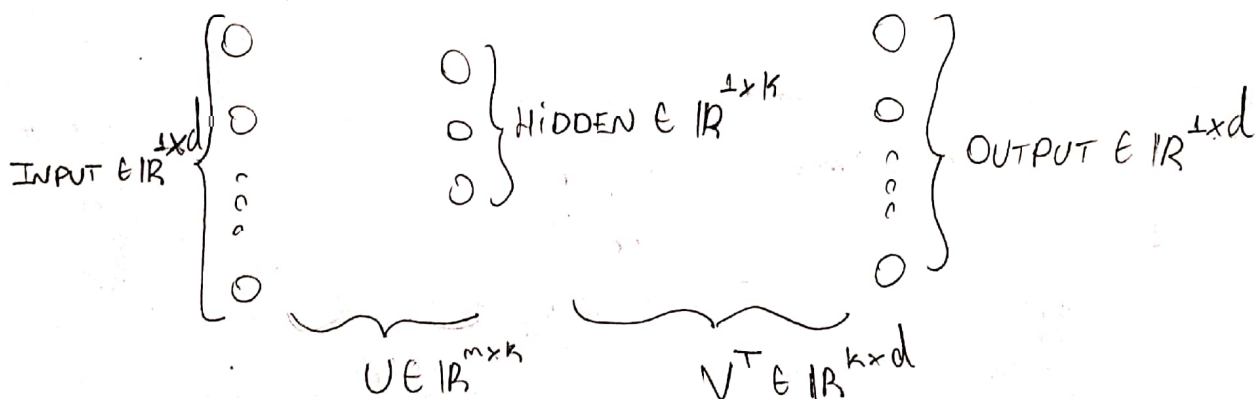
- ⑥ Apenas aprenderá a função identidade com alguns ruídos.
- ⑦ Então que o autoencoder aprenda a função identidade ou um mapeamento direto. O objetivo é aprender um espaço latente,
- ⑧ $\|x\|^2 = \sum_{i=1}^m (x_i)^2 \leq 1 \rightarrow$ normalização do espaço de entrada.

Para que o produto interno ou a semelhança dos vetores seja a maior possível, o vetor W também possuirá norma unitária.

- ⑨ Como o encoder é uma convolução, o decoder terá uma convolução transposta com filtro de mesmo tamanho que o filtro de entrada.

⑩ $D \in \mathbb{R}^{m \times d}$
 $U \in \mathbb{R}^{m \times k}$
 $V \in \mathbb{R}^{d \times k}$

$\rightarrow D \approx UV^T$



THIAGO MALTA LOUTINHO - 2014123335

① Sparse autoencoder: mapeia a entrada para um espaço de maior dimensão. Esse espaço é esparsos e o objetivo é gerar vetores base por meio da Loss com regularização L_1 .

Denoising Autoencoder: o objetivo é filtrar ruídos presentes na entrada dos dados. Dado um tipo de ruído que se quer filtrar, ele é adicionado a entrada do modelo e a saída é a entrada sem ruído. A rede então aprende a remover os ruídos.

Variational Autoencoder: O modelo aprende a distribuição dos dados de entrada a partir da modelagem de uma distribuição de variáveis latentes da entrada.

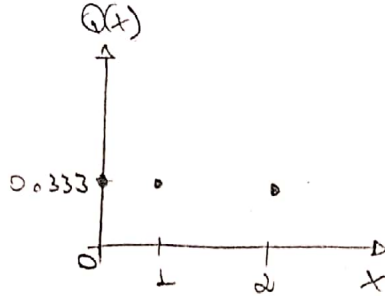
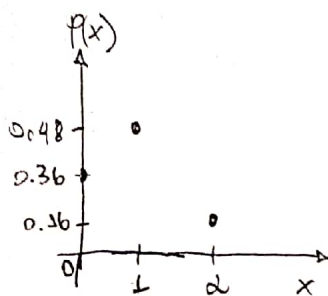
Contrastive Autoencoder: A função de perda é desmembrada para que pequenas variações na entrada produza pequenas variações no espaço latente.

② Um modelo generativo é um modelo que aprende a distribuição dos dados, possibilitando gerar novos exemplos da distribuição aprendida. O autoencoder é um modelo que mapeia uma entrada para uma saída através de um espaço latente, ele não aprende as distribuições dos dados. O autoencoder variacional aprende as distribuições do espaço latente e da saída, por isso é generativo.

$$\textcircled{3} \quad l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\phi(z|x_i)} \log[p_\theta(x_i|z)] + \kappa \mathcal{KL}(q_\phi(z|x_i) || p(z))$$

$$\textcircled{4} \quad D_{KL}(P||Q) = -\sum_x P \log Q + \sum_x P \log P$$

⑤



$$D_{KL}(Q||P) = \sum_x P_Q(x) \log \frac{P_Q(x)}{P_P(x)} = 0.333 \cdot \left[\log \left(\frac{0.333}{0.36} \right) + \log \left(\frac{0.333}{0.48} \right) + \log \left(\frac{0.333}{0.16} \right) \right]$$

$$D_{KL}(P||Q) = 0.36 \log \left(\frac{0.36}{0.333} \right) + 0.48 \log \left(\frac{0.48}{0.333} \right) + 0.16 \log \left(\frac{0.16}{0.333} \right)$$

$$D_{KL}(Q||P) \neq D_{KL}(P||Q)$$

⑦

$$\mu_\phi(x) = \text{ReLU}(2 \cdot \text{ReLU}(3x-1))$$

$$\Sigma_\phi(x) = \text{ReLU}(\text{ReLU}(3x+1)+1)$$

$$\mu_\theta(z) = \text{ReLU}(-z+12)$$

$$\Sigma_\theta(z) = \text{ReLU}(2z-10)$$

a) $x=2$

$$\epsilon = [0.5, 0, -0.5]$$

$$z_i = \mu_\phi(x) + \Sigma_\phi(x) \odot \epsilon_i$$

$$z_1 = \underbrace{\text{ReLU}(2 \cdot \text{ReLU}(6-1))}_{10} + \underbrace{\text{ReLU}(\text{ReLU}(6+1)+1)}_8 \cdot 0.5 = 10+4=14$$

$$z_2 = 10 + 8 \cdot 0 = 10$$

$$z_3 = 10 + 8 \cdot (-0.5) = 10-4=6$$

$$\vec{z} = [14, 10, 6]$$

b) $x_i = \mu_\theta(z), \Sigma_\theta(z)$

$x_2 = N(2, 10)$

$$x_1 = [\text{ReLU}(-z+12), \text{ReLU}(2z-10)] = [0, 18]$$

$x_3 = N(6, 2)$

$$x_1 = N(0, 18)$$

⑥ Separando a parte determinística da aleatória. A parte determinística é a média, o desvio padrão e a incerteza multiplicativa. A parte estocástica é modelada como um "ruído" normal com média zero e desvio padrão unitário.

⑧ $p(x)$ é intratável, pois não conseguimos estimá-lo diretamente.

⑨ ???