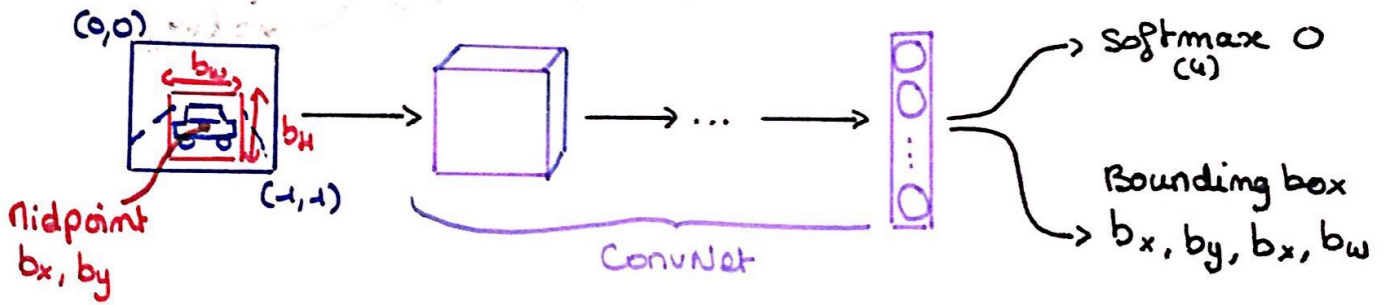


Détection d'objets (1)

=> L'algorithme est responsable de la classification avec localisation de l'image (i.e. mettre un cadre autour de l'objet) = 1 objet.

± Problème de détection = plusieurs objets à détecter dans l'image (de ± catégo).

Classification avec localisation



=> New target vector

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \left\{ \begin{array}{l} \text{probabilité qu'il y ait un objet qu'on} \\ \text{essaie de détecter (loss = logistic regression)} \\ \text{bounding box (loss = squared error)} \end{array} \right.$$

si $p_c = 1$, probabilité d'être la classe 1
2 ou 3... (loss = log like)

Si il n'y a pas d'objets, $p_c = 0$ et le reste vaut ? pour "don't care"

Loss

$$\mathcal{L}(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_s - y_s)^2 & \text{si } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{si } y_1 = 0 \end{cases}$$

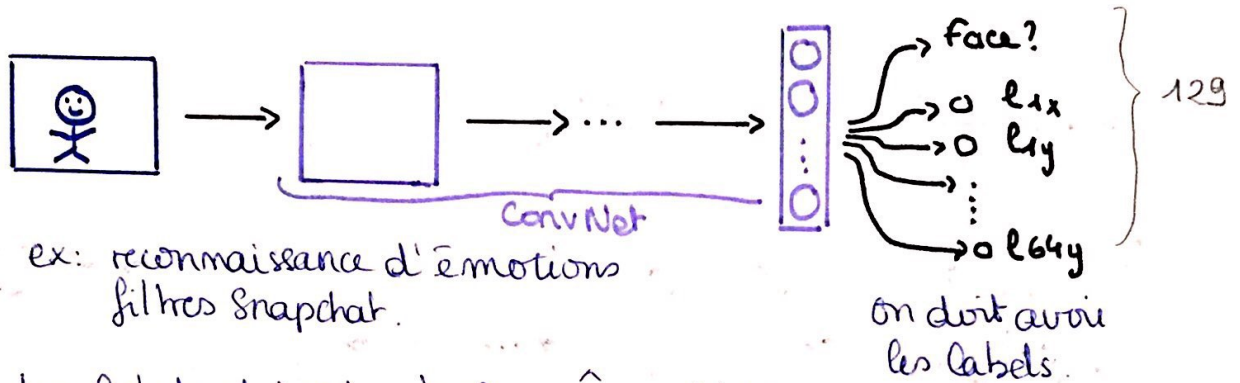
(on utilise l'erreur au carré pour simplifier).

↳ on s'intéresse qu'à la précision de l'algo pour estimer y_1

landmark detection

=> Algo qui reconnaît des points particuliers du visage (par exemple, les contours de l'œil, de la bouche).

Rend les coordonnées (x, y) des points.



ex: reconnaissance d'émotions
filtres Snapchat.

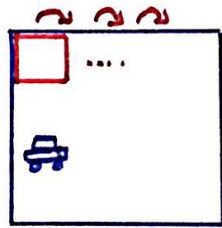
des labels doivent rester les mêmes d'une img à une autre.

Sliding windows detection

Idee:

(1) Entraîner un ConvNet à reconnaître une voiture en prédisant juste 0 ou 1 (c'est une voiture ou non)

(2)



choisir une boîte et pour chaque partie de l'image, dire si c'est une voiture ou pas en la faisant glisser.

On peut ensuite prendre une boîte plus large.

Paramètre à définir: le stride.

Inconvénients: coût de calcul.

Petit stride => grand nb de calcul => bonne perf.

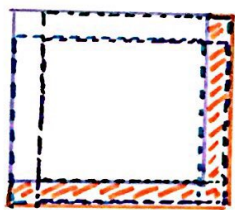
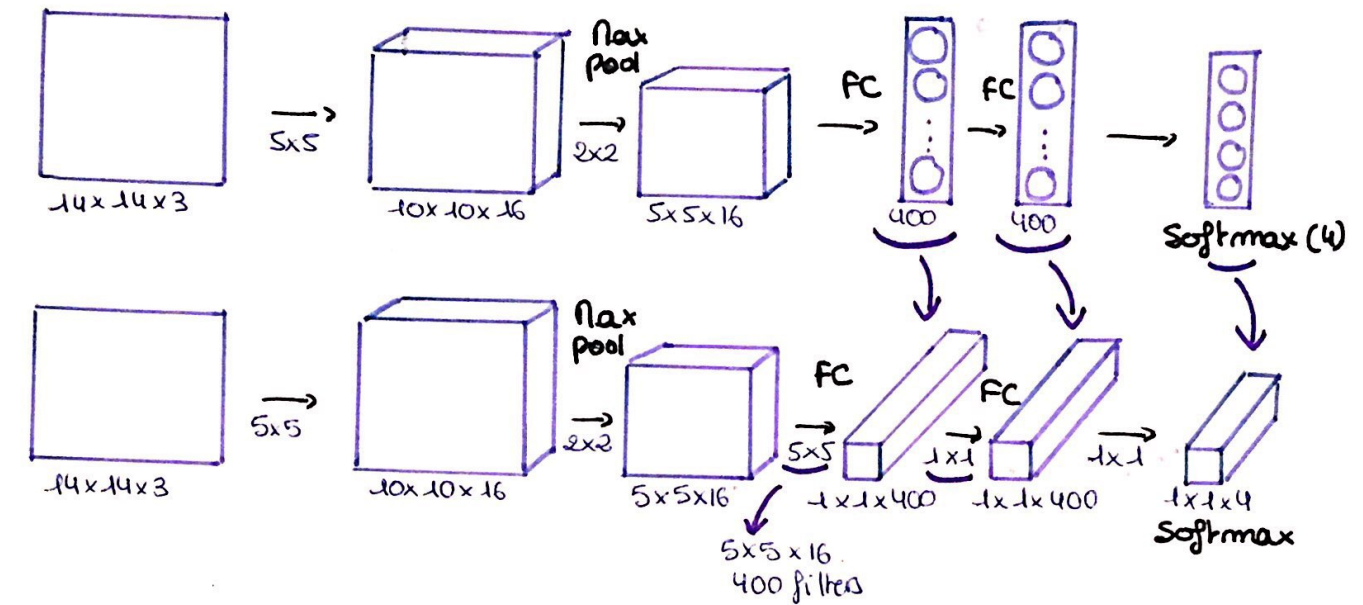
Grand stride => petit nb de calcul => peut v. perf.

Bonne méthode mais il existe des méthodes plus efficaces.

Détection d'objets (?)

Implementation de la fenêtre glissante avec convolution

=> Conversion de couches FC en couches de convolution.

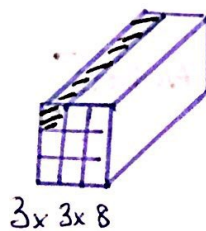
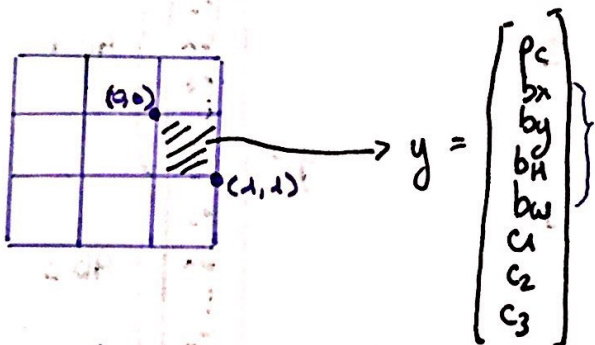


Convolution sur toute l'image

A la place de faire un Convnet sur chaque sous image, il suffit de faire un Convnet sur toute l'image car on obtient les résultats dans l'image finale.

Inconvénients : la position de la boîte ne sera pas précise.

=> Solution : **YOLO algorithm** (You only look once) = met une grille sur l'image et analyse chaque cellule de la grille.



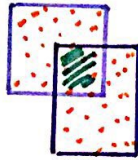
Target sortie

Le réseau de neurones => rend des "bounding box" précise

d'objet est assigné qu'à une seule des cellules.

b_x et b_w peuvent être > 1 si la boîte est sur plusieurs cellules.

Intersection over union



$$IoU = \frac{\text{Taille de l'intersection}}{\text{Taille de l'union}}$$

En général, si $IoU > 0,5 \Rightarrow$ "correct" bounding box
Plus l'IoU est élevé, plus la bounding box est précise.
On peut choisir un autre seuil

IoU = mesure de recouvrement entre deux bounding box
à quel point les bounding box sont similaires.

Non-max suppression

\Rightarrow permet de s'assurer que l'algo détecte chaque objet une seule fois.

(1) On supprime toutes les boîtes qui ont une $p_c \leq \text{seuil}$

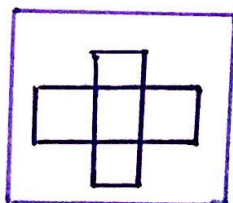
(2) Tant qu'il reste des boîtes :

- Sélectionner la boîte avec la plus gde p_c
- Retirer toutes les boîtes qui ont une $IoU \geq \text{seuil}$ avec la boîte précédente.

\Rightarrow Si plusieurs objets de \neq catégories, il faut faire tourner cet algo plusieurs fois.

Anchor boxes

\Rightarrow Pour que chaque cellule détecte plusieurs objets.



on peut les choisir à la main ou par k-means.

Anchor 1

Anchor 2



Chaque objet est assigné à la cellule qui contient son point "milieu" et l'anchor box qui a le plus gd IoU

$$y = \left[\begin{array}{c} p_c \\ b_x \\ b_y \\ \vdots \\ c_2 \\ c_3 \\ p_c \\ b_x \\ \vdots \\ c_2 \\ c_3 \end{array} \right] \left. \vphantom{\begin{array}{c} p_c \\ b_x \\ b_y \\ \vdots \\ c_2 \\ c_3 \\ p_c \\ b_x \\ \vdots \\ c_2 \\ c_3 \end{array}} \right\} \begin{array}{l} AB1 \\ \\ AB2 \end{array}$$

face recognition

face verification VS

face recognition

- Input image, name/ID
- Sortie : l'image correspond à la personne

1:1

- Base de données de k personnes
- Sortie : ID si l'image est parmi les k personnes

1:K

=> Construire un système de face verification et si la précision est suffisamment grande, l'utiliser dans un système de face recognition.

One-shot learning

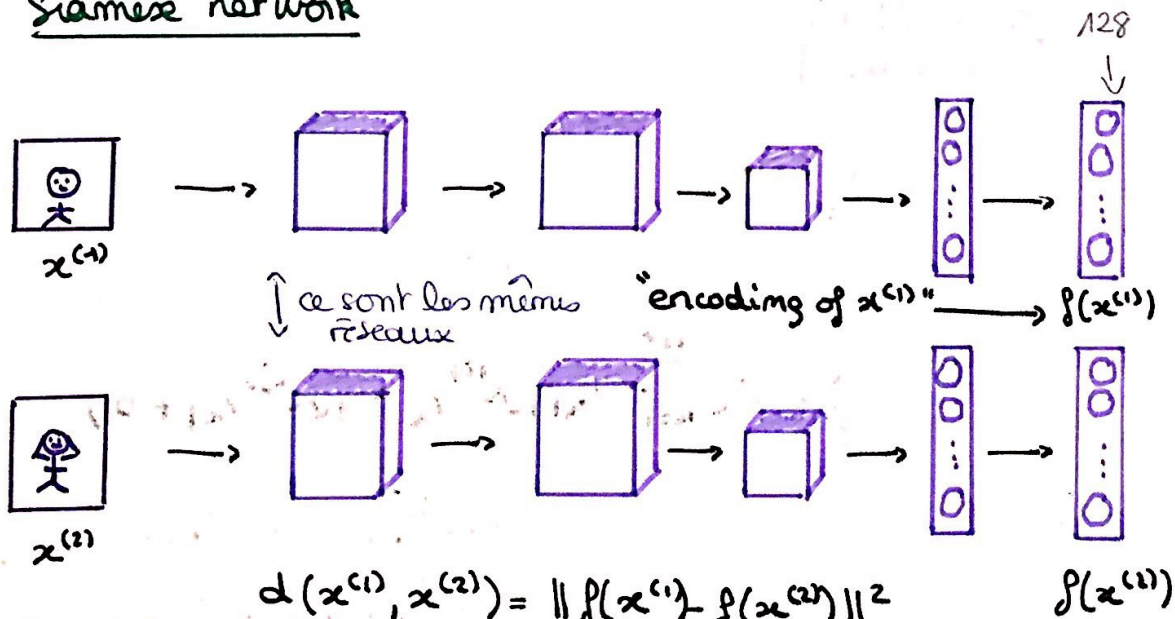
=> Pour construire un syst de reconnaissance faciale, on doit être capable de reconnaître une personne avec seulement 1 image.

Si on entraîne le CONVNET avec des images des \neq personnes et qu'on utilise une softmax en sortie, ça ne marchera pas très bien car petit jeu de données ou alors à réentraîner à chq fois qu'on a une nouvelle personne.

=> on va plutôt apprendre une fonction de similarité

$d(\text{img1}, \text{img2}) = \text{degré de } \neq \text{ entre les images}$

Siamese network



la distance est petite si il s'agit de la même personne.

Triplet Loss

Etant donné 3 images $\overbrace{A, P, N}$

A: anchor
P: positive
N: Negative

$$\mathcal{L}(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

permet de
garder la loss
positive

Hyperparamètre α
permet d'éloigner
la paire positive de la
paire négative
= margin

Fonction de perte :

$$J = \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)})$$

=> on a besoin d'un jeu de données avec plusieurs paires d'images de la même personne.

Si on choisit les triplets aléatoirement => trop de triplets seront trop faciles pour l'algo

=> la descente de gradient ne fera pratiquement rien

On choisit des triplets difficiles

↗ l'efficacité
du calcul

=> l'algo va vraiment essayer d'éloigner les deux paires en faisant en sorte qu'il y ait une marge entre les deux.

Binary classification

=> Pour prédire si les deux images sont les mêmes (=1)
ou pas (=0)

dalle
du vecteur
de sortie
(en exemple)

$$\hat{y} = \sigma\left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b\right)$$

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k} = \chi^2$$