# FCN: (Baseline)

$$\varepsilon(\Theta) = \sum e(X_\Theta(p), \ell(p)) \rightarrow \text{Objective function}$$

$p$ = pixel index, $\ell(p)$ = gt label, $X_\Theta(p)$ = net labelin

$e(\ell(p), X_\Theta(p))$ = per pixel loss

$\Theta$ = network parametrization; updated w/ SGD and
backprop

# BoxSup:
### Overlapping

Objective Function $\mathcal{E}_O = \frac{1}{N} \sum_S (1 - IoU(B,S)) \delta(\ell_B, \ell_S)$   ①   → Find S w/ greatest overlap w/ B and $\ell_B = \ell_S$

$S$ = candidate segment mask

$B$ = gt bounding box annotation

$IoU(B,S) \in [0,1] \rightarrow$ intersection-over-union ratio

$\uparrow IoU \Rightarrow \uparrow$ box-candidate mask overlap

$$\delta(\ell_B, \ell_S) = \begin{cases} 1 & \text{if } \ell_B = \ell_S \\ 0 & \text{otherwise} \end{cases} \qquad \begin{cases} \ell_B = \text{semantic label of bounding box } B \\ \ell_S = \text{semantic label of candidate segment } S \end{cases}$$

Minimizing $\mathcal{E}_O$ implies higher IoUs for consistent semantic labels

$N$ = # of candidate segments

---

$\mathcal{E}_r = \sum_p e(X_\theta(p), \ell_S(p))$ ,   ②

$\ell_S(p)$ = semantic label at pixel $p$ used for network training

Target of regression: estimated candidate segment

---

Overarching Objective Function: $\mathcal{E} = \min_{\theta, \{\ell_S\}} \sum_i (\mathcal{E}_O + \lambda \mathcal{E}_r)$   ③

$\sum_i$ = sum over all images

$\lambda = 3$ (fixed weighting parameter)

Parameters to optimize: a) net parameters $\theta$

                b) labelling of all candidate segments $\{\ell_S\}$

## Full Supervision Loss Function:

$I$ = set of pixels of image; $N$ = # of pixels

$S_{ic}$ = CNN score for pixel $i$ and class $c$

Softmax probability of $c$ at $i$: $S_{ic} = \dfrac{e^{S_{ic}}}{\sum\limits_{k=1}^{N} e^{S_{ik}}} \in [0,1]$

$G$ = ground truth map

↳ pixel $i$ belongs to class $G_i$

Loss on single training image:

Cross-entropy loss ①

$$L^{pix}(S,G) = -\sum_{i \in I} \log (S_{iG_i})$$

(if $G_i$ undefined, set $\log(S_{iG_i}) = 0$ for that value of $i$)

## Image-Level Supervision Loss Function:

$\{1,\dots,N\}$ = set of all classes CNN trained to recognize

$L \subseteq \{1,\dots,N\}$ ~~known~~ classes present in image

$L' \subseteq \{1,\dots,N\}$ classes not present in image

$$L_{img}(S,L,L') = -\frac{1}{|L|}\sum_{c \in L} \log(S_{t_c c}) - \frac{1}{|L''|}\sum_{c \in L'} \log(1 - S_{t_c c}) \quad ②$$

, where $t_c = \underset{i \in I}{\arg\max}\, S_{ic}$

↳ Single-image cross-entropy loss

## Point-Level Supervision Loss Function:

Combines ① and ②

↑ (① only for supervised points)

$I_s$ = set of pixels w/ known class; supervised pixels

$$L_{point}(S,G,L,L') = L_{img}(S,L,L') - \sum_{i \in I_s} a_i \log(S_{iG_i})$$

$a_i$ = relative importance of each supervised pixel

# Point-level Supervision w/ Object Prior:

$P_i$ = probability pixel $i$ belongs to an object

$O$ = set of object classes; $O'$ = set of background classes

e.g. PASCAL VOC $\Rightarrow O$ = set of 20 object classes

$O'$ = generic background class

$$L_{obj}(S, P) = -\frac{1}{|I|} \sum_{i \in I} \left[ P_i \log\left(\sum_{c \in O} S_{ic}\right) + (1 - P_i) \log\left(1 - \sum_{c \in O} S_{ic}\right)\right]$$