



Department of
Computer Science

香港城市大學
City University of Hong Kong



ICLR
International Conference On
Learning Representations



哈爾濱工業大學(深圳)
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN
计算机科学与技术学院

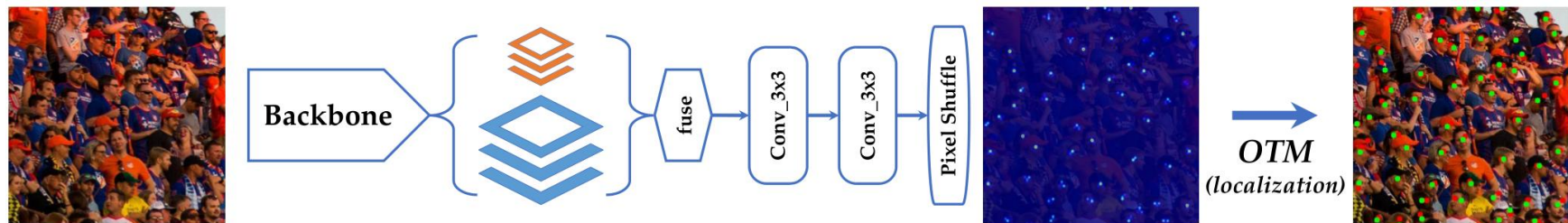
Proximal Mapping Loss: Understanding Loss Functions in Crowd Counting & Localization

Wei Lin¹, Jia Wan², and Antoni B. Chan¹

¹Department of Computer Science, City University of Hong Kong

²School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen
elonlin24@gmail.com, jiawan1998@gmail.com, abchan@cityu.edu.hk

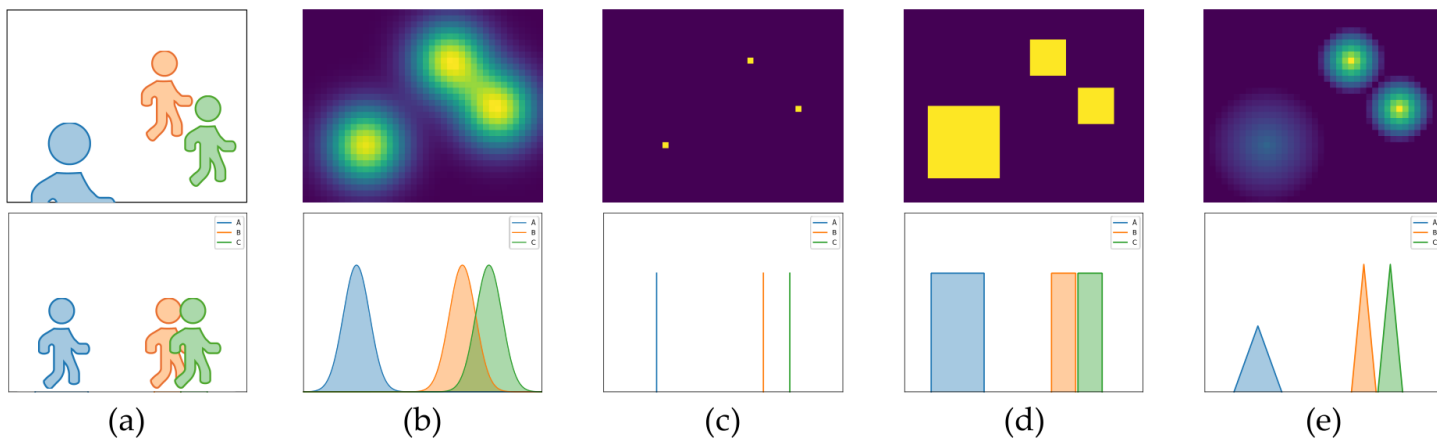
Supervised Crowd Counting



Supervised crowd counting is normally formulated as a regression task:

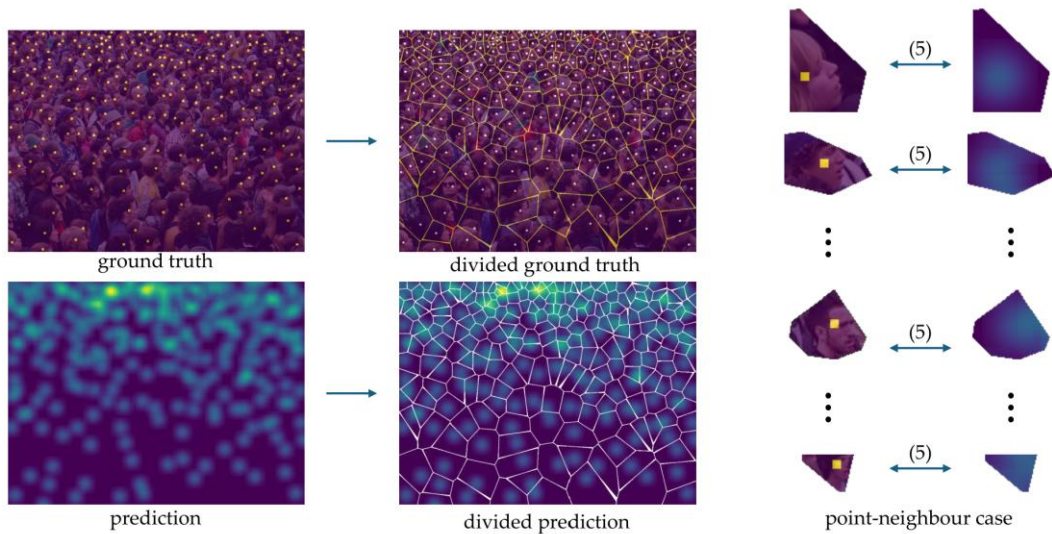
- **Input:** an image contain crowds;
 - **Output:** a density map demonstrate the distribution and count of crowd in the input;
 - **Ground Truth:** a point map in which a pixel with a value of one denotes a person's location.
-
- **OTM** is applied to regress the density map into a point map.
 - **Applications:** video surveillance and public safety services, traffic congestion control, marine environmental monitoring

Supervised Crowd Counting



- (a) A synthetic input with three humans;
- (b) Density regression with the intersection hypothesis (using a Gaussian prior), where one pixel may correspond to multiple objects;
- (c) Point prediction;
- (d) Head region segmentation;
- (e) Density regression without the intersection hypothesis (proposed method), where one pixel corresponds to one object.

Proximal Mapping Loss

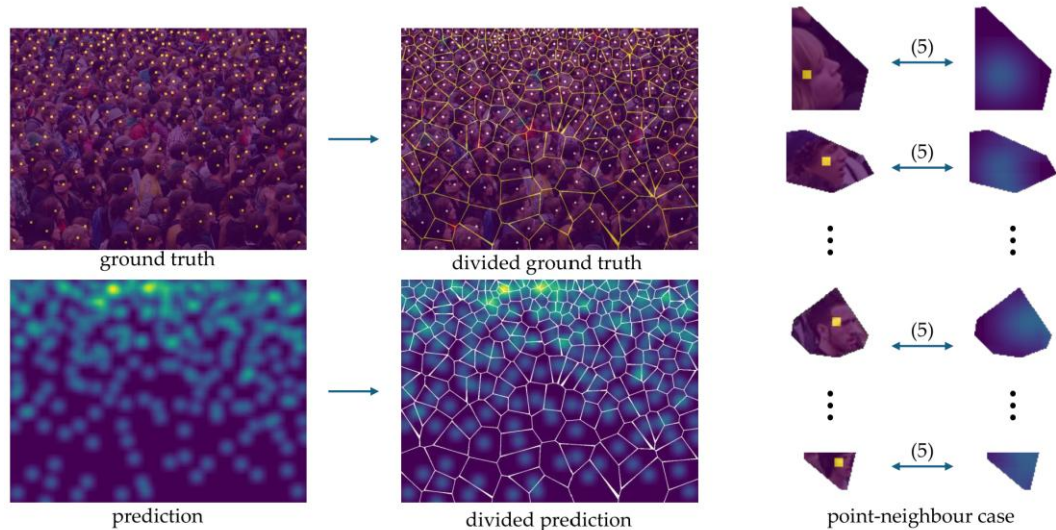


(I) Divide Stage: assign each pixel to its nearest GT point

(II) Conquer Stage: loss computation

- nearest neighbour is adopted to split the density map into multiple irregular patches without overlapping.
- In PML, the loss computation is divided into multiple simpler sub-problems, as each point-neighbour case can be handled independently.

Proximal Mapping Loss



(I) **Divide Stage:** assign each pixel to its nearest GT point

(II) **Conquer Stage:** loss computation

$$\mathcal{L}(\mathcal{A}, \mathcal{B}) = \sum_{j=1}^m \tilde{\mathcal{L}}(\tilde{\mathcal{A}}_j, \mathbf{b}_j), \quad \tilde{\mathcal{A}}_j = \{(a_i, \mathbf{x}_i)\}_{i \in \mathcal{X}_j} \quad \mathbf{b}_j = (1, \mathbf{y}_j), \quad (1)$$

$$\mathcal{X}_j = \{i \mid \|\mathbf{x}_i - \mathbf{y}_j\|_2 \leq \|\mathbf{x}_i - \mathbf{y}_k\|_2, \quad \forall \mathbf{y}_k \in \mathcal{B}\}, \quad (2)$$

Proximal Mapping Loss

By defining $\tilde{\mathbf{a}} = [\tilde{a}_i]_i^{\tilde{n}}$ constructed from $\tilde{\mathcal{A}}$, the objective inherited from GL is to minimize the transport $f(\tilde{\mathbf{a}}) = \mathbf{c}^\top \tilde{\mathbf{a}}$, where $\mathbf{c} = [c_i]_{i=1}^{\tilde{n}}$ measures the cost when moving a unit mass from \mathbf{x}_i to \mathbf{y} .

Proximal mapping: $\tilde{\mathbf{a}}_{t+1} \approx \operatorname{argmin}_{\mathbf{p}} \underbrace{f(\tilde{\mathbf{a}}_t) + \nabla f(\tilde{\mathbf{a}}_t)^\top (\mathbf{p} - \tilde{\mathbf{a}}_t)}_{\text{linear approximation of } f(\tilde{\mathbf{a}}_{t+1})} + \underbrace{\frac{\tau}{2} \|\mathbf{p} - \tilde{\mathbf{a}}_t\|^2}_{\text{regularizer}}$

$$\left. \begin{array}{l} \nabla f(\tilde{\mathbf{a}}_t) = \mathbf{c} \\ \xi \subseteq \{\mathbf{p} \mid \mathbf{p}^\top \mathbf{1} = 1, \mathbf{p} \in \mathbb{R}^{\tilde{n}}\} \end{array} \right\} \Rightarrow \mathcal{L}(\tilde{\mathcal{A}}, \mathbf{b}) = \min_{\mathbf{p} \in \xi} \mathbf{c}^\top \mathbf{p} + \frac{\tau}{2} \|\mathbf{p} - \tilde{\mathbf{a}}\|^2$$

$$\text{Bregman divergence } \mathcal{D}_\varphi(\mathbf{p}, \tilde{\mathbf{a}}) \Rightarrow \mathcal{L}(\tilde{\mathcal{A}}, \mathbf{b}) = \min_{\mathbf{p} \in \xi} \mathbf{p}^\top \mathbf{a} + \tau \mathcal{D}_\varphi(\mathbf{p}, \tilde{\mathbf{a}})$$

Proximal Mapping Loss

$$\mathcal{L}(\tilde{\mathcal{A}}, \mathbf{b}) = \min_{\mathbf{p} \in \xi} \mathbf{p}^\top \mathbf{a} + \tau \mathcal{D}_\varphi(\mathbf{p}, \tilde{\mathbf{a}})$$

loss function	τ	$\mathcal{D}_\varphi(\mathbf{p}, \mathbf{a})$	ξ	\mathbf{p}^*
L2 loss (Zhang et al., 2016)	0	$\frac{1}{2} \ \mathbf{p} - \mathbf{a}\ ^2$	$\mathcal{N}(\mu \Sigma)$	$\mathcal{N}(0 \sigma \mathbf{1}_{2 \times 2})$
Bayesian loss (Ma et al., 2019)	$\frac{1}{ \mathbf{1}^\top \mathbf{a} - 1 }$	$\frac{1}{2} \ \mathbf{p} - \mathbf{a}\ ^2$	-	$\mathbf{a} - \frac{1}{ \mathbf{1}^\top \mathbf{a} - 1 } \mathbf{c} + \eta$
P2PNet (Song et al., 2021)	-	$\ \mathbf{p} - \mathbf{a}\ _1$	$\delta(\cdot)$	$\delta(\arg \min_j \mathbf{c}_j - \tau \mathbf{a}_j)$
DM-Count (Wang et al., 2020a)	∞	$\text{KL}(\mathbf{p} \mid \mathbf{a})$	-	$\mathbf{a} / \ \mathbf{a}\ _1$

PML & L2 loss

$$\mathcal{L}(\tilde{\mathcal{A}}, \mathbf{b}) = \min_{\mathbf{p} \in \xi} \mathbf{p}^\top \mathbf{a} + \tau \mathcal{D}_\varphi(\mathbf{p}, \tilde{\mathbf{a}}) \quad \text{with } \mathcal{D}_\varphi = \frac{1}{2} \|\mathbf{p} - \tilde{\mathbf{a}}\|^2$$



$$\mathcal{L}_2 = \min_{\mathbf{p}} \mathbf{c}^\top \mathbf{p} + \frac{\tau}{2} \|\mathbf{p} - \tilde{\mathbf{a}}\|_2^2, \quad s.t. \quad \mathbf{p}^\top \mathbf{1} = \sum_{i=0}^{\tilde{n}} p_i = 1.$$



$$\mathbf{p}^* = \tilde{\mathbf{a}} - \frac{1}{\tau_1} \mathbf{c} + \eta, \quad \eta = \frac{1}{\tilde{n}} \left[1 - \left(\tilde{\mathbf{a}} - \frac{1}{\tau} \mathbf{c} \right)^\top \mathbf{1} \right]$$

Here η takes the role as “filler” such that $\mathbf{p}^{*\top} \mathbf{1} = 1$.



$$\frac{\partial \mathcal{L}_2}{\partial \tilde{a}_i} = c_i - \tau \eta \quad \Rightarrow \quad \mathcal{L}_2 = \frac{\tau}{2} \|\mathbf{a} - \text{detach}(\mathbf{p}^*)\|_2^2$$

PML & L2 loss

$$\mathcal{L}_2 = \frac{\tau}{2} \|a - \text{detach}(\mathbf{p}^*)\|_2^2 \quad \text{where} \quad \mathbf{p}^* = a - \frac{1}{\tau} \mathbf{c} + \eta$$

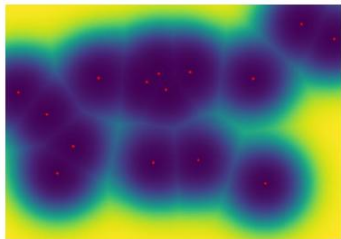
$$\text{Dynamic L2 loss} \quad \begin{cases} \lim_{\tau \rightarrow 0} \mathbf{p}^* = \delta(\mathbf{y}) \\ \lim_{\tau \rightarrow \infty} \mathbf{p}^* = \tilde{\mathbf{a}} + (1 - \mathbf{1}^\top \tilde{\mathbf{a}}) \end{cases}$$

$$\text{Traditional L2 loss} \quad \begin{cases} \tau = 0 \\ \xi \subseteq \mathcal{N}(\boldsymbol{\mu} | \Sigma) \quad \Rightarrow \quad \mathbf{p}^{*'} \leftarrow \mathcal{N}(\boldsymbol{\mu} \mid (x - \boldsymbol{\mu})^\top (x - \boldsymbol{\mu})), \quad \boldsymbol{\mu} = x^\top \mathbf{p}^* \\ \Sigma \succcurlyeq \sigma^2 \mathbf{I}_{2 \times 2} \quad \Rightarrow \quad \mathbf{p}^{*'} \leftarrow \mathcal{N}(\mathbf{y}_j | \sigma^2 \mathbf{I}_{2 \times 2}) \end{cases}$$

PML & Bayesian Loss



(a) image & density map



(b) Bayesian loss

$$\mathcal{L}_b = \sum_{i=1}^n q(y_0|x_i)a_i + \sum_{j=1}^m \left| \sum_{i=1}^n q(y_j|x_i)a_i - 1 \right|$$

→

$$\mathcal{L}_b = \underbrace{\mathbf{q}^\top \mathbf{a}}_{\text{background loss}} + \underbrace{\left| \mathbf{a}^\top \mathbf{1} - 1 \right|}_{\text{count loss}},$$

- **count loss** forces the sum of \mathbf{a} to be close to 1;
- **background loss** forces the distribution of \mathbf{a} to be close to $\delta(\mathbf{y})$;

$$q_{[i]} = q(y_0|x_i) = \frac{q(x_i|y_0)}{q(x_i|y_0)+1}, \quad q(x_i|y_0) \stackrel{\text{def.}}{=} \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(d-c_i)^2}{2\sigma^2} \right]$$

$$q_{[i]} \propto q(x_i|y_0) \quad \Rightarrow \quad q_{[i]} > q_{[j]} \text{ if } c_i > c_j$$

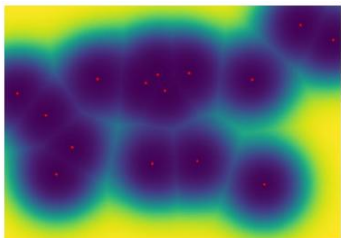
PML & Bayesian Loss

$$\mathcal{L}(\tilde{\mathcal{A}}, b) = \min_{p \in \xi} p^\top a + \tau \mathcal{D}_\varphi(p, \tilde{a}) \quad \text{with} \quad \mathcal{D}_\varphi = \frac{1}{2} \|p - \tilde{a}\|^2$$

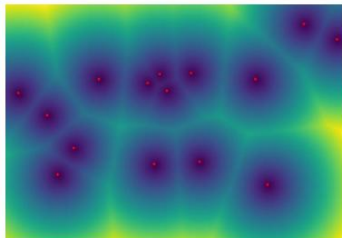
$$\Rightarrow p^* = \tilde{a} - \frac{1}{\tau_1} c + \eta, \quad \eta = \frac{1}{\tilde{n}} \left[1 - \left(\tilde{a} - \frac{1}{\tau} c \right)^\top \mathbf{1} \right]$$



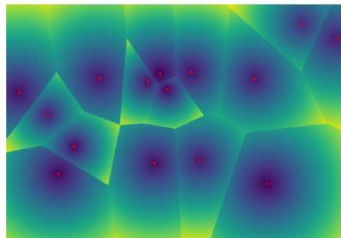
(a) image & density map



(b) Bayesian loss



(c) generalized loss



(d) proximal mapping loss

$$\frac{\partial \mathcal{L}_2}{\partial \tilde{a}_i} = c_i - \tau \eta = c_i - \frac{1}{\tilde{n}} c^\top \mathbf{1} + \frac{\tau}{\tilde{n}} (\tilde{a}^\top \mathbf{1} - 1)$$

$$\Rightarrow \mathcal{L}_2 = \underbrace{(c - \bar{c})^\top a}_{\text{background loss}} + \frac{\tau}{2\tilde{n}} \underbrace{(\tilde{a}^\top \mathbf{1} - 1)^2}_{\text{count loss}},$$

PML & Bayesian Loss

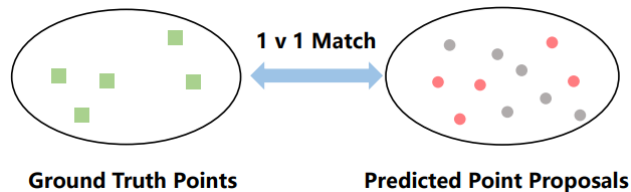
$$\mathcal{L}_2 = \underbrace{(\mathbf{c} - \bar{\mathbf{c}})^\top \mathbf{a}}_{\text{background loss}} + \frac{\tau}{2\tilde{n}} \underbrace{\left(\tilde{\mathbf{a}}^\top \mathbf{1} - 1 \right)^2}_{\text{count loss}},$$

$$\tau = \text{detach} \left(2 / |\mathbf{1}^\top \tilde{\mathbf{a}} - 1| \right) \begin{cases} \mathcal{L}'_b = (\mathbf{c} - \bar{\mathbf{c}}) \mathbf{a} + \frac{1}{\tilde{n}} |\mathbf{1}^\top \tilde{\mathbf{a}} - 1| \\ \mathbf{p}^* = \mathbf{a} - \left(\frac{1}{2} |\mathbf{1}^\top \tilde{\mathbf{a}} - 1| \right) \mathbf{c} + \eta \end{cases}$$

L1 norm is robust to noise annotation:

- If the predicted count is close to 1, \mathbf{p}^* will be close to the distribution of \mathbf{a} ;
- If the count is far from GT, \mathbf{p}^* will be close to the distribution of $\delta(\mathbf{y})$

PML & P2PNet



The matching is implemented via Hungarian algorithm with the cost matrix:

$$\mathcal{D}(\mathcal{P}, \hat{\mathcal{P}}) = (\tau \|p_i - \hat{p}_j\|_2 - \hat{c}_j)_{i \in N, j \in M}$$

$$\mathcal{L}_{p2p} = \min_{\mathbf{p} \in \xi_{p2p}} \mathbf{c}^\top \mathbf{p} + \tau \|\mathbf{p} - \tilde{\mathbf{a}}\|_1, \quad \xi_{p2p} = \{\delta(x_i) | x_i \in \tilde{\mathcal{A}}\}$$



$$\begin{aligned} \mathcal{L}_{p2p} &= \min_{\mathbf{p} \in \xi_{p2p}} \mathbf{c}^\top \mathbf{p} + \tau [\mathbf{p}^\top (1 - \mathbf{a}) + (1 - \mathbf{p})^\top \tilde{\mathbf{a}}] \\ &= \min_{\mathbf{p} \in \xi_{p2p}} \underbrace{(\mathbf{c} - 2\tau \tilde{\mathbf{a}})^\top \mathbf{p}}_{\text{matching strategy}} + \underbrace{\tau (\mathbf{p}^\top \mathbf{1} + \tilde{\mathbf{a}}^\top \mathbf{1})}_{\text{constant } \tau(1 + \|\tilde{\mathbf{a}}\|_1)}, \end{aligned}$$



$$\mathbf{p}^* = \delta(\operatorname{argmin}_j c_j - 2\tau \tilde{a}_j)$$

PML & DMC

$$\mathcal{L}_{dmc} = \underbrace{\mathbf{c}^\top \frac{\mathbf{a}}{\|\mathbf{a}\|_1}}_{\text{OT loss}} + \tau \underbrace{|\mathbf{1}^\top \mathbf{a} - 1|}_{\text{count loss}}$$

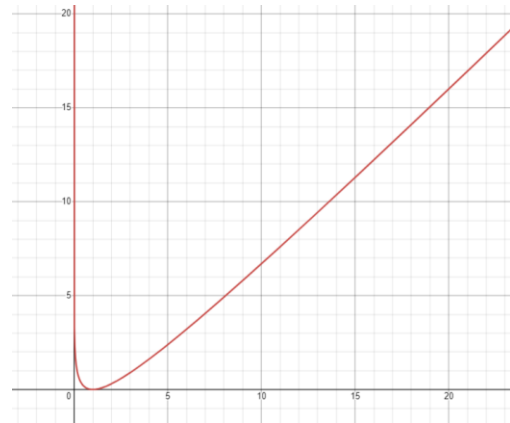
$$\varphi(\mathbf{x}) = \mathbf{x}^\top \log \mathbf{x} - \mathbf{x}^\top \mathbf{1} \quad \Rightarrow \quad \mathcal{D}_\varphi(\mathbf{p}, \mathbf{a}) = \mathbf{p}^\top \log \frac{\mathbf{p}}{\mathbf{a}} - \mathbf{p}^\top \mathbf{1} + \mathbf{a}^\top \mathbf{1}.$$

$$\Rightarrow \mathcal{L}_{dmc} = \min_{\mathbf{p}} \mathbf{c}^\top \mathbf{p} + \tau (\mathbf{p}^\top \log \frac{\mathbf{p}}{\mathbf{a}} - \mathbf{p}^\top \mathbf{1} + \mathbf{a}^\top \mathbf{1}) \quad \text{s.t.} \quad \mathbf{p}^\top \mathbf{1} = 1,$$

$$\Rightarrow p_i^* = \frac{a_i \exp(-c_i/\tau)}{\eta}, \quad \eta = \sum_j a_j \exp(-c_j/\tau).$$

η also serves as “filler”, ensuring the sum of elements in \mathbf{p}^* equals 1.

$$\tau \rightarrow \infty \quad \Rightarrow \quad \mathbf{p} = \frac{\mathbf{a}}{\|\mathbf{a}\|_1} \quad \Rightarrow \quad \mathcal{L}_{dmc}^{(\varphi)} = \underbrace{\mathbf{c}^\top \frac{\mathbf{a}}{\|\mathbf{a}\|_1}}_{\text{OT loss}} + \tau \underbrace{(\mathbf{1}^\top \mathbf{a} - \log(\mathbf{1}^\top \mathbf{a}) - 1)}_{\text{count loss}}$$



Experiments

METHOD	(backbone)	ShTech A		ShTech B		UCF-QNRF		JHU ++		NWPU	
		MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
MCNN (Zhang et al., 2016)		232.5	714.6	110.2	173.2	277.0	426.0	188.9	483.4	232.5	714.6
CSRNet (Li et al., 2018)	(VGG-16)	68.2	115.0	10.6	16.0	110.6	190.1	85.9	309.2	121.3	387.8
SFCN (Wang et al., 2019)	(ResNet-101)	64.8	107.5	7.6	13.0	102.0	171.4	77.5	297.6	105.7	424.1
BL (Ma et al., 2019)	(VGG-19)	62.8	101.8	7.7	12.7	88.7	154.8	75.0	299.9	105.4	454.2
KDMG (Wan et al., 2020)	(VGG-19)	63.8	99.2	7.8	12.7	99.5	173.0	69.7	268.3	100.5	415.5
DMC (Wang et al., 2020a)	(VGG-19)	59.7	95.7	7.4	11.8	85.6	148.3	68.4	283.3	88.4	357.6
NoiseCC (Wan & Chan, 2020)	(VGG-19)	61.9	99.6	7.4	11.3	85.8	150.6	67.7	258.5	96.9	534.2
P2PNet (Song et al., 2021)	(VGG-16bn)	52.7	85.6	6.3	9.9	85.3	154.5	-	-	77.4	362.0
UOTCC (Ma et al., 2021)	(VGG-19)	58.1	95.9	6.5	10.2	83.3	142.3	60.5	252.7	87.8	387.5
GL (Wan et al., 2021)	(VGG-19)	61.3	95.4	7.3	11.7	84.3	147.5	59.9	259.5	79.3	346.1
ChfL (Shu et al., 2022)	(VGG-19bn)	57.5	94.3	6.9	11.0	80.3	137.6	57.0	235.7	76.8	343.0
PET (Liu et al., 2023)	(VGG-16bn)	49.3	78.8	6.2	9.7	79.5	144.3	58.5	238.0	74.4	328.5
STEERER (Han et al., 2023)	(HRNet)	54.5	86.9	<u>5.8</u>	<u>8.5</u>	<u>74.3</u>	<u>128.3</u>	<u>54.3</u>	238.1	63.7	<u>309.3</u>
PML (ours)	(VGG-16bn)	<u>50.6</u>	<u>80.7</u>	6.1	9.7	79.5	142.7	58.9	249.6	75.7	353.1
PML (ours)	(VGG-19)	55.5	89.0	6.0	9.3	76.6	132.2	57.4	227.4	73.6	338.6
PML (ours)	(HRNet)	52.3	84.7	5.4	8.2	73.2	127.5	52.6	<u>230.8</u>	<u>63.8</u>	306.9

Table 2: Comparison of our PML with recent crowd counting methods.

Experiments

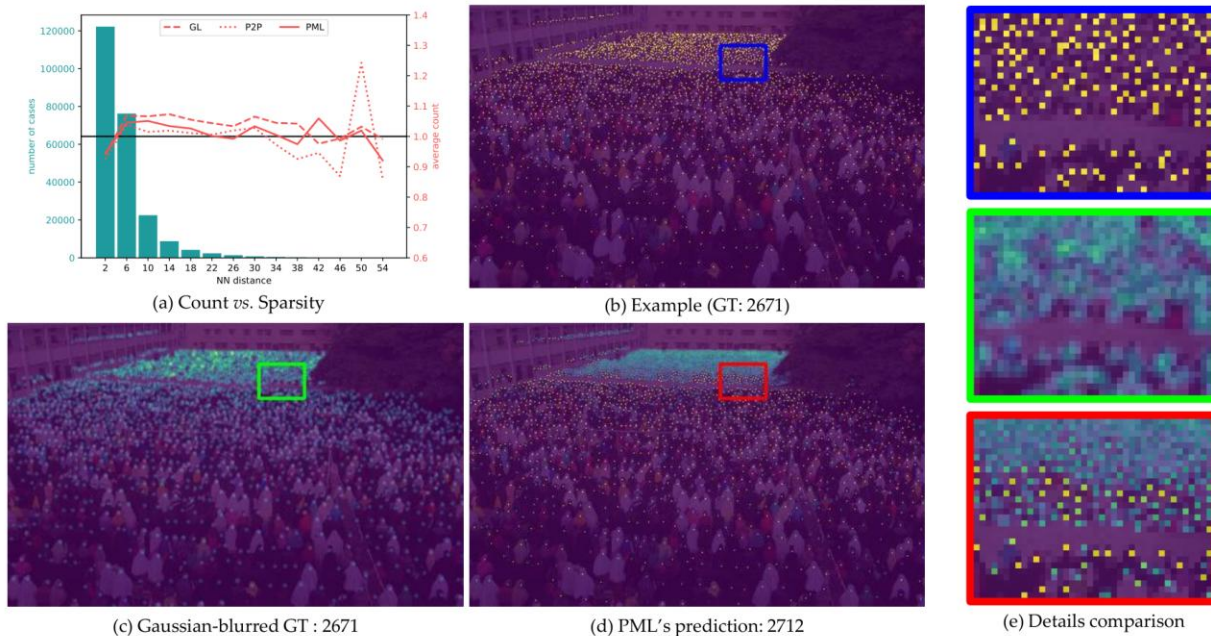


Figure 7: The performance difference of PML when handling sparse and dense crowds.

Experiments

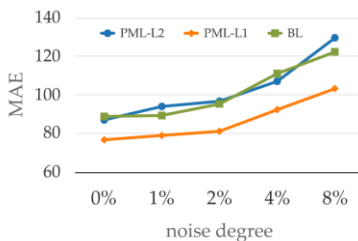
	F1-meas.	Prec.	Rec.
RAZNet	0.599	0.666	0.543
GL+LM	0.660	0.800	0.562
GL+OTM	0.683	0.710	0.658
P2PNet	0.729	0.676	0.685
PET	0.742	0.752	0.732
STEERER+LM	0.770	0.814	0.730
PML(VGG-19)+OTM	0.735	0.776	0.698
PML(HRNet)+OTM	0.802	<u>0.809</u>	0.795
PML(HRNet)+LM	<u>0.790</u>	0.803	<u>0.777</u>

Table 3: Localization on NWPU-Crowd.

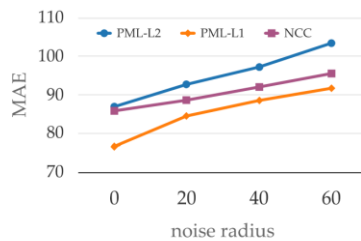
	MCNN		CSRNet		VGG19		HRNet	
	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
L2 Loss	186.4	283.6	110.6	190.1	98.7	176.1	92.03	157.49
BL	190.6	272.3	107.5	184.3	88.8	154.8	85.52	149.55
NoiseCC	177.4	259.0	96.5	163.3	85.8	150.6	-	-
DMC	176.1	263.3	103.6	180.6	85.6	148.3	82.07	144.84
GL	142.8	227.9	92.0	165.7	84.3	147.5	78.37	140.23
GCFL	-	-	83.0	139.8	80.3	137.6	-	-
PML (ours)	138.9	215.7	82.1	139.0	77.6	132.8	73.17	127.45

Table 4: Comparison of loss functions and backbones on UCF-QNRF dataset.

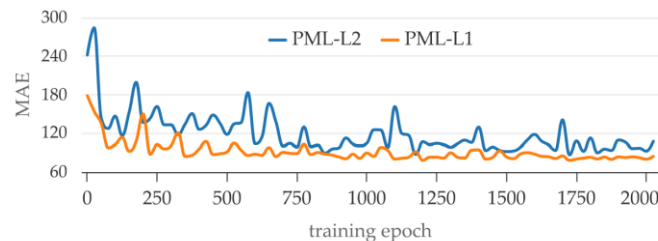
Experiments



(a) Adding noise following BL



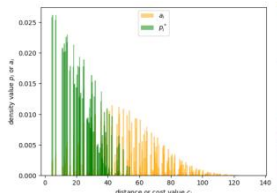
(b) Adding noise following NCC



(c) MAE of PML with L2- and L1-norm during training



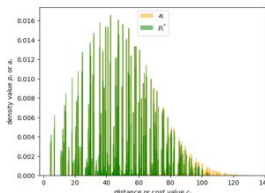
(a) Image & prediction



(c) Under estimation (L1-norm)



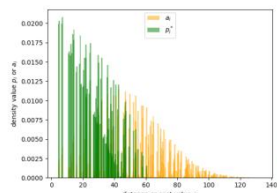
(d) Well estimation (L1-norm)



(e) Over estimation (L1-norm)



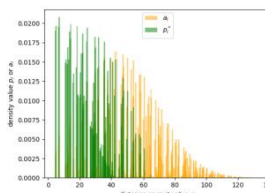
(b) GT & neighbor pixels



(f) Under estimation (L2-norm)



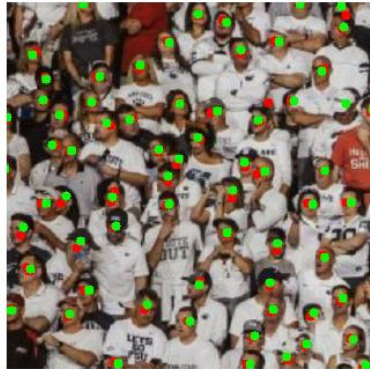
(g) Well estimation (L2-norm)



(h) Over estimation (L2-norm)



Experiments





Department of
Computer Science

香港城市大學
City University of Hong Kong



ICLR
International Conference On
Learning Representations



哈爾濱工業大學(深圳)
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN
计算机科学与技术学院

Thanks

Wei Lin¹, Jia Wan², and Antoni B. Chan¹

¹Department of Computer Science, City University of Hong Kong

²School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen

elonlin24@gmail.com, jiawan1998@gmail.com, abchan@cityu.edu.hk