

Learning Human Poses from Actions

Aditya Arun¹, C.V. Jawahar¹, M. Pawan Kumar^{2,3}



INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY

HYDERABAD

¹International Institute of Information Technology, Hyderabad ²University of Oxford, ³The Alan Turing Institute

The **Alan Turing** Institute

JNIVERSITY OF

DXFORI

Aim

- We propose a novel framework to **learn human poses** using **diverse data set**.
- In our setting, diverse data set includes:
 - A few expensively annotated images with ground-truth joint annotations
 - Several cheaply annotated images specifying human action.

Actions and Poses



Learning Objective

• The diversity coefficient is the expected value of the loss.

$$\operatorname{PIV}_{\Delta}(\operatorname{Pr}_{\theta}, \operatorname{Pr}_{\mathbf{w}}) = \sum_{\mathbf{h}^{k}, \mathbf{h}^{k'} \in \mathcal{H}} \Delta(\mathbf{h}^{k}, \mathbf{h}^{k'}) \operatorname{Pr}_{\theta}(\mathbf{h}^{k}) \operatorname{Pr}_{\mathbf{w}}(\mathbf{h}^{k'})$$

Diversity of a distribution Pr $\mathrm{DIV}_{\Delta}(\mathrm{Pr}_{\theta}, \mathrm{Pr}_{\mathbf{w}}) = \Delta(\mathbf{h}^1, \mathbf{h}^2) \operatorname{Pr}_{\theta}(\mathbf{h}^1) \operatorname{Pr}_{\theta}(\mathbf{h}^2)$

Pr₆(h¹|x,a)





Fig: Images from MPII Human Pose test set for four different action classes.

- Human poses for each action class are sufficiently different.
- Probabilistic formulation necessary due to high intra-class variance.

Probabilistic Formulation

Tasks:

- During training, model the uncertainty in the pose for every action.
- During inference, predict the pose given an image.

Two separate distributions for two tasks^[1]:

. A **Conditional Distribution** of the pose given an image and its corresponding action.

h¹' Diversity between the distributions Pr_e and Pr_w: Pr_e(h¹|x,a) $\mathrm{DIV}_{\Delta}(\mathrm{Pr}_{\theta}, \mathrm{Pr}_{\mathbf{w}}) = \Delta(\mathbf{h}^1, \mathbf{h}^{2'}) \mathrm{Pr}_{\theta}(\mathbf{h}^1) \mathrm{Pr}_{\mathbf{w}}(\mathbf{h}^{2'})$

- We design a joint learning objective that minimizes the dissimilarity coefficient^[4] (DISCO) between the prediction distribution and the conditional distribution.
- Formally, the dissimilarity coefficient is written as an affine combination of:
 - \circ diversity between Pr_{\bullet} and Pr_{w} ; and
 - diversity of each distribution.

 $\mathrm{DISC}_{\Delta}(\mathrm{Pr}_{\mathbf{w}}, \mathrm{Pr}_{\boldsymbol{\theta}}) = \mathrm{DIV}_{\Delta}(\mathrm{Pr}_{\mathbf{w}}, \mathrm{Pr}_{\boldsymbol{\theta}}) - \gamma \mathrm{DIV}_{\Delta}(\mathrm{Pr}_{\mathbf{w}}, \mathrm{Pr}_{\mathbf{w}})$ $-(1-\gamma)\mathrm{DIV}_{\Delta}(\mathrm{Pr}_{\theta},\mathrm{Pr}_{\theta})$

Optimization

Pr_e(h²|x,a)

We estimate the parameters of the two networks in three stages:

- Supervised training of the two networks using small amount of ground truth pose data.
- Iterative training: updating one network while keeping the other fixed on diverse data.
- Jointly optimize both the networks together on diverse data.

Visualization of the effect of iterative learning of the two networks are shown below:



Prediction Distribution: A distribution over the pose given an image.

Ideally, the two distributions should match exactly.

Pr_w(h|x)

We use a probabilistic network, DISCO Nets^[2], to model parameters of these distributions.

DISCO Stacked Hourglass Network

By adding noise filter to the stacked hourglass network^[3], we construct the DISCO net.



















Iter 0 Iter 2 Iter 5 Iter 7 Fig: Example of superimposed pose predictions by (a) prediction network and (b) conditional network illustrating the uncertainty in the pose across training iterations.

Experiments and Results

• Diverse data set:

- We use MPII Human Pose data set.
- To obtain the various diverse data set, we choose three different data splits, {25-75,50-50,75-25}%, where we randomly discard 75%, 50%, and 25% of the pose annotations from the training images respectively.

Fig: For a single input image x and three different noise samples $\{z^1, z^2, z^3\}$ (represented as red, green and blue matrix respectively), the network produces three different candidate poses $\{h^1, h^2, h^3\}$.

Pointwise Prediction:

- Single input **x**.
- Multiple $\{\mathbf{h}^1, \mathbf{h}^2, \dots, \mathbf{h}^k\}$ sampled from the probabilistic hourglass network.
- Optimal prediction according to the Δ with maximum expected utility (EU):

$$\mathbf{h}_{\Delta}^{*}(\mathbf{x};\mathbf{w}) = \arg\max_{k \in [1,K]} EU(\mathbf{h}^{k}) = \arg\min_{k \in [1,K]} \sum_{k'=1}^{K} \Delta(\mathbf{h}^{k},\mathbf{h}^{k'})$$

References

- 1. M. Pawan Kumar, Ben Packer, and Daphne Koller. *Modeling latent variable uncertainty for loss-based* learning. In ICML, 2012.
- Diane Bouchacourt, M. Pawan Kumar, and Sebastian Nowozin. DiscoNets: Dissimilarity coefficient 2. networks. In NIPS, 2016.
- 3. Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In ECCV, 2016.
- C. Radhakrishna Rao. Diversity and dissimilarity coefficients: a unified approach. Theoretical 4. population biology, 21(1):24–43, 1982.

• Methods:

• Fully supervised (**FS**) stacked hourglass network trained on supervised subset.

- Non probabilistic pointwise network (**PW**) trained with diverse data set. Ο
- Probabilistic DISCO-HG network (**Pr**) trained with diverse data set.

Results:

