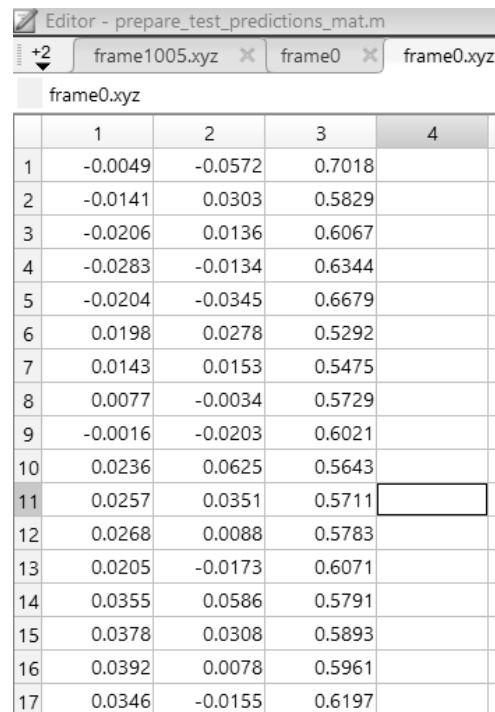


Normalization of RGB hand dataset RHD

Overview of dataset

- Dataset split
 1. Train: 41258
 2. Test: 2728
- Each image is composed of
 1. Left hand (21 joints)
 2. Right hand (21 joints)
- Ground Truth Annotation
 1. 3D (unit: probably *m*?)
 2. Camera Parameters
 3. 2D projection
 4. Visible mask

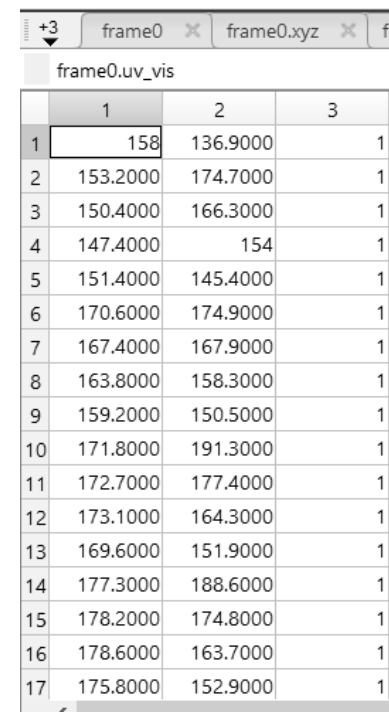


Editor - prepare_test_predictions_mat.m

+2 frame1005.xyz x frame0 x frame0.xyz

frame0.xyz

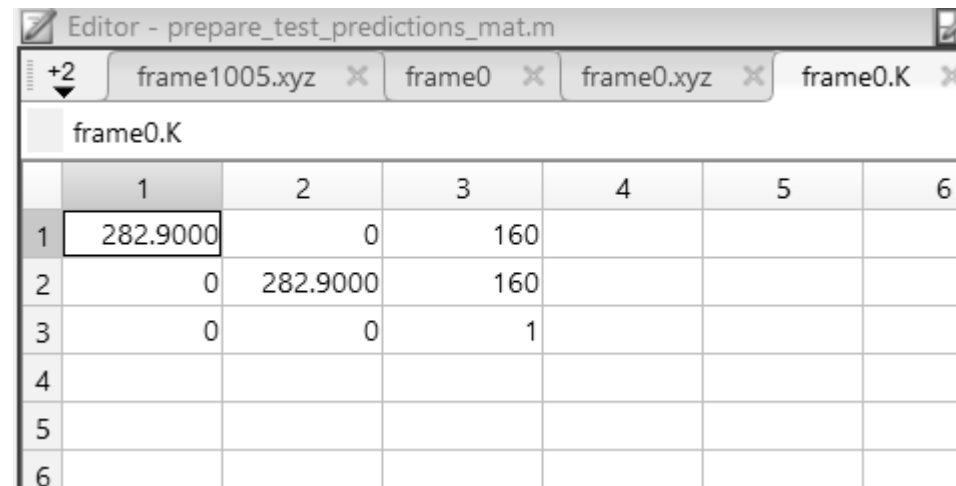
	1	2	3	4
1	-0.0049	-0.0572	0.7018	
2	-0.0141	0.0303	0.5829	
3	-0.0206	0.0136	0.6067	
4	-0.0283	-0.0134	0.6344	
5	-0.0204	-0.0345	0.6679	
6	0.0198	0.0278	0.5292	
7	0.0143	0.0153	0.5475	
8	0.0077	-0.0034	0.5729	
9	-0.0016	-0.0203	0.6021	
10	0.0236	0.0625	0.5643	
11	0.0257	0.0351	0.5711	
12	0.0268	0.0088	0.5783	
13	0.0205	-0.0173	0.6071	
14	0.0355	0.0586	0.5791	
15	0.0378	0.0308	0.5893	
16	0.0392	0.0078	0.5961	
17	0.0346	-0.0155	0.6197	



+3 frame0 x frame0.xyz x fr

frame0.uv_vis

	1	2	3
1	158	136.9000	1
2	153.2000	174.7000	1
3	150.4000	166.3000	1
4	147.4000	154	1
5	151.4000	145.4000	1
6	170.6000	174.9000	1
7	167.4000	167.9000	1
8	163.8000	158.3000	1
9	159.2000	150.5000	1
10	171.8000	191.3000	1
11	172.7000	177.4000	1
12	173.1000	164.3000	1
13	169.6000	151.9000	1
14	177.3000	188.6000	1
15	178.2000	174.8000	1
16	178.6000	163.7000	1
17	175.8000	152.9000	1



Editor - prepare_test_predictions_mat.m

+2 frame1005.xyz x frame0 x frame0.xyz x frame0.K x

frame0.K

	1	2	3	4	5	6
1	282.9000	0	160			
2	0	282.9000	160			
3	0	0	1			
4						
5						
6						

Overview of dataset

- Camera Parameters

$$K = \begin{bmatrix} f_x & 0 & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

Where f_x is focus length (x direction), f_y is focus length (y direction)

In the dataset, as far as I am concerned, f_x is equal f_y

u_0 is x direction offset of raw image, and v_0 is y direction offset

$$KP = proj$$

$$\begin{bmatrix} f_x & 0 & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} \frac{x_k}{z_k} \\ \frac{y_k}{z_k} \\ 1 \end{bmatrix} = \begin{bmatrix} f_x \frac{x_k}{z_k} + u_0 \\ f_y \frac{y_k}{z_k} + v_0 \end{bmatrix} = \begin{bmatrix} u_k \\ v_k \end{bmatrix} \text{ (for instance, joint } k \text{)}$$

2.5D Pose Representation

- \times Predict absolute 3D in *camera coordinate* is infeasible given the projection and scale ambiguity.
- *Instead*, predict 2.5D pose

$$P_k^{2.5D} = (u_k, v_k, Z_k^r)$$

Where u_k, v_k are image coordinates

Z_k^r is root-relative depth value (relative location to root joint)

2.5D Pose Representation

- One thing to mention is that different hand has different scales
Remove the scale ambiguity by

Scale Normalization (P is 3D before scale normalization)

$$\hat{P} = \frac{C}{s} P$$

C is a constant set to 1.0, $s = \|P_n - P_{parent(n)}\|_2$ is the bone length of a specific bone (palm root -> index MCP here)

Thus the bone length after normalization is

$$\hat{s} = \|\hat{P}_n - \widehat{P_{parent(n)}}\|_2 \text{ which should be the constant } C$$

2.5D Pose Representation

- Look at the formula

$$KP = proj$$
$$K\hat{P} = \widehat{proj} = proj$$

The scale-normalized 2.5D pose we want to predict is

$$P_k^{2.5D} = (\widehat{u}_k, \widehat{v}_k, \widehat{Z}_k^r) = (u_k, v_k, \widehat{Z}_k^r)$$

About u_k, v_k (projection on raw image)

- the network predicts local projection normalized in $[0, 1]$, which can be recovered back to global location provided with bounding box crop information

2.5D Pose Representation

About \widehat{Z}_k^r

- The paper says that all hand centers are approximately in a range between 40 cm and 65 cm,
- The real values of each joint's depth lies approximately in a range 0.40 to 0.65,
- Thus one can analyze the training dataset to crop a rough cube, for each sample, which contains the hand propitiously (the cube size in optical axis direction **should be a constant value**)
- And then normalize the “scale-normalized root-relative” depth in that range

Overall Architecture

Infer 2.5D from
CNN

- Scale-normalized root-relative joint prediction get

Scale normalized
root Z

- Depth value of scale-normalized root $\widehat{Z_{root}}$ get

Scale recovery

- Recover the scale-normalized absolute joint prediction to absolute w/o scale norm.

Scale-normalized root-relative joint

- Input: cropped monocular RGB hand image I
- Output: scale-normalized root-relative joint in 2.5D space

$$\widehat{P^{2.5D}} = \left\{ (u_k, v_k, \widehat{Z}_k^r) \right\}_{k \in K}$$

Methodology: CNN

$$f(I) = \widehat{P^{2.5D}}$$

- ✓ Regression
- ✓ Heatmap (2D/3D)
- ✓ Other map representation (fully conv network)

Scale-normalized depth value of root

- Estimate \widehat{Z}_{root}

Let n be index MCP, m be parent of n : palm root

If intrinsic camera parameter is unknown, multiple 3D solutions can have same 2D projection

Given camera parameters f_x, f_y, u_0, v_0 , one unique solution:

$$(\widehat{X}_n - \widehat{X}_m)^2 + (\widehat{Y}_n - \widehat{Y}_m)^2 + (\widehat{Z}_n - \widehat{Z}_m)^2 = c^2$$

(remember the aforementioned **scale normalization** step)

Scale-normalized depth value of root

$$\bullet \begin{pmatrix} f_x & 0 & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} \frac{x_k}{z_k} \\ \frac{y_k}{z_k} \\ 1 \end{pmatrix} = \begin{pmatrix} f_x \frac{x_k}{z_k} + u_0 \\ f_y \frac{y_k}{z_k} + v_0 \end{pmatrix} = \begin{pmatrix} u_k \\ v_k \end{pmatrix}$$

$$\widehat{x}_n = \frac{(u_n - u_0)\widehat{Z}_n^r}{f_x}, \widehat{x}_m = \frac{(u_m - u_0)\widehat{Z}_m^r}{f_x}, \widehat{y}_n = \frac{(v_n - v_0)\widehat{Z}_n^r}{f_y}, \widehat{y}_m = \frac{(v_m - v_0)\widehat{Z}_m^r}{f_y}$$

$$\text{Let } \widehat{Z}_{root} = t$$

Scale-normalized depth value of root

$$\begin{aligned} & \bullet \left(\frac{u_n - u_0}{f_x} (t + \widehat{z}_n^r) - \frac{u_m - u_0}{f_x} (t + \widehat{z}_m^r) \right)^2 + (\widehat{X}_n - \widehat{X}_m)^2 \\ & \bullet \left(\frac{v_n - v_0}{f_y} (t + \widehat{z}_n^r) - \frac{v_m - v_0}{f_y} (t + \widehat{z}_m^r) \right)^2 + (\widehat{Y}_n - \widehat{Y}_m)^2 \\ & \bullet \left((t + \widehat{z}_n^r) - (t + \widehat{z}_m^r) \right)^2 - C^2 = 0 \quad (\widehat{Z}_n - \widehat{Z}_m)^2 - C^2 \end{aligned}$$

Which means

Scale-normalized depth value of root

- $\left(\left(\frac{u_n - u_m}{f_x} \right) t + \frac{u_n - u_0}{f_x} \widehat{Z}_n^r - \frac{u_m - u_0}{f_x} \widehat{Z}_m^r \right)^2 +$
- $\left(\left(\frac{v_n - v_m}{f_y} \right) t + \frac{v_n - v_0}{f_y} \widehat{Z}_n^r - \frac{v_m - v_0}{f_y} \widehat{Z}_m^r \right)^2 +$
- $\left(\widehat{Z}_n^r - \widehat{Z}_m^r \right)^2 - C^2 = 0$

To quadratic formula

Scale-normalized depth value of root

$$\bullet \left(\frac{u_n - u_m}{f_x} \right)^2 t^2 + \left(\frac{v_n - v_m}{f_y} \right)^2 t^2 + 2t \left(\left(\frac{u_n - u_m}{f_x} \right) \left(\frac{u_n - u_0}{f_x} \widehat{Z}_n^r - \right.$$

Scale-normalized depth value of root

- Let r_A be $\left(\frac{u_n - u_m}{f_x}\right)^2 + \left(\frac{v_n - v_m}{f_y}\right)^2$

$$r_B \text{ be } 2\left(\left(\frac{u_n - u_m}{f_x}\right)\left(\frac{u_n - u_0}{f_x}\widehat{Z}_n^r - \frac{u_m - u_0}{f_x}\widehat{Z}_n^r\right) + \left(\frac{v_n - v_m}{f_y}\right)\left(\frac{v_n - v_0}{f_y}\widehat{Z}_n^r - \frac{v_m - v_0}{f_y}\widehat{Z}_n^r\right)\right)$$

Scale-normalized depth value of root

- Then $r_A t^2 + r_B t + r_C = 0$

$$t = \frac{-r_B \pm \sqrt{r_B^2 - 4r_A r_C}}{2r_A}$$

Only one t is valid (*bigger one*)

Scale-normalized depth value of root

- Screenshot of code
Validated on several training samples using the
 1. Ground truth 2D projection
 2. Ground truth 3D (for scale norm calculation)
 3. Camera parameters

To solve the scale-normalized depth value of root \widehat{Z}_{root}

Scale-normalized depth value of root

```
0 0
0 -0.00384871
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 7.230918 7.139670
Root 2 vs Ground truth scale normalized root: -6.859144 7.139670
-----
0 1
1 -0.0693331
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 14.902864 13.141487
Root 2 vs Ground truth scale normalized root: -10.641583 13.141487
-----
1 0
2 0.0976509
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 3.929582 4.751347
Root 2 vs Ground truth scale normalized root: -6.891019 4.751347
-----
1 1
3 0.332327
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 8.937823 10.315970
Root 2 vs Ground truth scale normalized root: -8.927805 10.315970
-----
2 0
4 -0.00321121
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 5.424101 5.402974
Root 2 vs Ground truth scale normalized root: -4.650224 5.402974
-----
2 1
3 0
6 -0.199604
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 6.308317 5.315545
Root 2 vs Ground truth scale normalized root: -6.527187 5.315545
-----
```

```
10 0.185482
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 3.717680 4.064990
Root 2 vs Ground truth scale normalized root: -3.632873 4.064990
-----
5 1
6 0
12 0.145443
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 3.481955 3.731007
Root 2 vs Ground truth scale normalized root: -3.408043 3.731007
-----
6 1
7 0
7 1
15 0.00237677
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 6.934124 6.949356
Root 2 vs Ground truth scale normalized root: -5.341604 6.949356
-----
8 0
16 -0.288994
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 7.704199 6.486199
Root 2 vs Ground truth scale normalized root: -7.700637 6.486199
-----
8 1
17 -0.112387
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 1.087305 0.984722
Root 2 vs Ground truth scale normalized root: 0.976353 0.984722
-----
9 0
18 -0.00622031
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root: 11.489633 11.163298
Root 2 vs Ground truth scale normalized root: 3.041926 11.163298
-----
```

Scale recovery

- Need to know the global hand scale

The big idea is that when estimated scale normalized root-relative 3D is scaled, the difference from mean bone length (averaged over all training samples) is minimized. (For each bone)

Formula:

$$\hat{s} = \underset{s}{\operatorname{argmin}} \sum_{k,l \in \mathcal{E}} \left(s \cdot \|\widehat{P}_k - \widehat{P}_l\|_2 - \mu_{kl} \right)^2$$

Where μ_{kl} is average bone length between keypoint k and keypoint l on training set.

Scale recovery

$$\begin{aligned} & \bullet \sum_{k,l \in \mathcal{E}} \left(scale \|\widehat{P}_k - \widehat{P}_l\|_2 - \mu_{kl} \right)^2 = \\ & \quad scale^2 \sum_{k,l \in \mathcal{E}} \left(\|\widehat{P}_k - \widehat{P}_l\|_2 \right)^2 + scale \cdot (-2) \cdot \sum_{k,l \in \mathcal{E}} (\|\widehat{P}_k - \widehat{P}_l\|_2 \cdot \mu_{kl}) + \\ & \quad \sum_{k,l \in \mathcal{E}} (\mu_{kl})^2 = \end{aligned}$$

$$s_A \, scale^2 + s_B \, scale + s_C$$

The equation reaches its nadir at $scale = \frac{s_B}{-2 \, s_A}$

Scale recovery

- Where

$$s_A = \sum_{k,l \in \mathcal{E}} \left(\|\widehat{P}_k - \widehat{P}_l\|_2 \right)^2,$$

$$s_B = (-2) \cdot \sum_{k,l \in \mathcal{E}} (\|\widehat{P}_k - \widehat{P}_l\|_2 \cdot \mu_{kl}),$$

$$s_C = \sum_{k,l \in \mathcal{E}} (\mu_{kl})^2$$

Quite elegant!

Scale recovery

```
Root 2 vs Ground truth scale normalized root:  -4.669377    4.840677
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.110433    0.126036
-----
10    1
21 -0.112877
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    5.978505    5.392392
Root 2 vs Ground truth scale normalized root:    -4.817154    5.392392
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.111744    0.123711
-----
11    0
22 -0.122755
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    10.824253    9.879951
Root 2 vs Ground truth scale normalized root:    -7.451607    9.879951
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.108607    0.090770
-----
11    1
23 -0.0621625
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    7.240659    6.911001
Root 2 vs Ground truth scale normalized root:    -8.195185    6.911001
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.110399    0.091926
-----
12    0
24 -0.0702172
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    6.068675    5.628127
Root 2 vs Ground truth scale normalized root:    -6.175155    5.628127
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.110438    0.093797
```

```
Solving global hand scale which is
Solved global scale vs real global scale:    0.111290    0.126024
-----
19    1
20    0
40 0.126782
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    6.728361    7.524141
Root 2 vs Ground truth scale normalized root:    -4.892658    7.524141
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.110965    0.090761
-----
20    1
21    0
42 -0.0113637
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    5.400281    5.359130
Root 2 vs Ground truth scale normalized root:    -5.337444    5.359130
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.104483    0.093709
-----
21    1
43 0.0769343
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    6.960660    7.445234
Root 2 vs Ground truth scale normalized root:    -2.686331    7.445234
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.104623    0.093724
-----
22    0
44 -0.100865
---
Solving scale normalized zroot
Root 1 vs Ground truth scale normalized root:    9.155630    8.706736
Root 2 vs Ground truth scale normalized root:    -6.145841    8.706736
-----
Solving global hand scale which is
Solved global scale vs real global scale:    0.110954    0.125765
```

Final step

- Suppose that CNN $f(I) = \{(\widehat{u}_k, \widehat{v}_k, \widetilde{\widehat{Z}_k^r})\}_{k \in K}$ takes as input the cropped image I , and predicts the **“normalized”** scale normalized root-relative 2.5D pose (**“normalized”** means normalized for CNN training: notation $\widetilde{}$)
 - local scale-norm projection (ranges in $[0, 1]$) $\widetilde{\widehat{u}_k} \widetilde{\widehat{v}_k}$
 - normed root-relative scale-norm depth (ranges in, for example, $[-1, 1]$) $\widetilde{\widehat{Z}_k^r}$

Final step

- First off, take an unnormalization effort to convert CNN output to real scale-normalized root-relative 2.5D pose (for example, $(\widehat{u}_k, \widehat{v}_k, \widehat{Z}_k^r) = (480 \text{ pixel}, 220 \text{ pixel}, 14\text{cm})$)

$$g(I) = \{(\widehat{u}_k, \widehat{v}_k, \widehat{Z}_k^r)\}_{k \in K} = \text{unnorm}(f(I)) = \text{unnorm}(\{(\widehat{u}_k, \widehat{v}_k, \widehat{Z}_k^r)\}_{k \in K})$$

- Afterwards, compute scale-normalized root depth \widehat{Z}_{root} ,
add each root-relative depth prediction \widehat{Z}_k^r with \widehat{Z}_{root}

$$\widehat{Z}_k = \widehat{Z}_k^r + \widehat{Z}_{root}$$

Final step

- Get scale normalized absolute $\widehat{X}_k, \widehat{Y}_k$

$$\widehat{X}_k = \frac{(\widehat{u}_k - u_0)\widehat{Z}_k^r}{f_x}, \widehat{Y}_k = \frac{(\widehat{v}_k - v_0)\widehat{Z}_k^r}{f_y}$$

As you can see, \widehat{X}_k and \widehat{Y}_k are computed using 2.5D \widehat{u}_k and \widehat{v}_k and \widehat{Z}_k^r
 \widehat{Z}_k is computed using \widehat{Z}_k^r and previously got \widehat{Z}_{root}

- To this end, we have acquired scale-normalized absolute 3D joint in camera frame (not root-relative any more)

$$\widehat{P} = \{\widehat{P}_k\}_{k \in K} = \{(\widehat{X}_k, \widehat{Y}_k, \widehat{Z}_k)\}_{k \in K}$$

Final step

- Use the statistics on average bone length (each bone) $\{\mu_{kl}\}_{k,l \in \mathcal{E}}$,

And scale-normalized absolute 3D pose \hat{P}

We solve the hand scale “*scale*”

- Finally, we come to the solution

$$P = \{P_k\}_{k \in K} = \left\{ \frac{scale}{C} \hat{P}_k \right\}_{k \in K}$$

- This is the global 3D in camera frame,
camera matrix \times *global 3D* = *global projection in raw iamge*