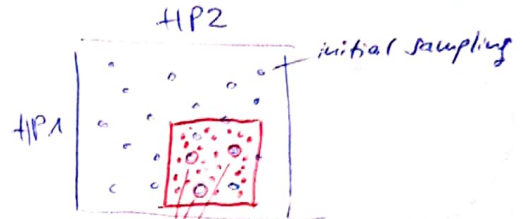


# Hyperparameter tuning

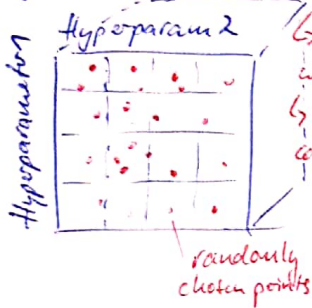
- 1st  $\alpha$  ← most important
- 2nd  $\beta = \sim 0.9$  should be ok
- 3rd  $\beta_1, \beta_2, \epsilon \rightarrow 0.9, 0.999, 10^{-8}$  — usually not varied much
- 2nd # layers
- 3rd # hidden units
- 2nd learning rate decay
- 2nd mini-batch size

Coarse to fine sampling



those worked best, is randomly sample around them → add finer sample points

Try random values : Don't use a grid!



Reasons: don't know up front which parameter is most important  
 will end up with more different configurations for  $\alpha$  than with a grid  
 ⇒ more richly exploring!  
 (like GoE)

Using an appropriate scale to pick hyperparameters (HP)

reasonable picking at random → not true for all HP!

→  $n = 50, \dots, 100 \rightarrow$

→ #layers  $L : 2-4 \rightarrow 2, 3, 4$

example  $\alpha : 0.0001, \dots, 1 \rightarrow$

rather go for log-scale

$\alpha = b \cdot \log_{10} 0.0001$   
 $a = -4$

Python:  $r = -4 \cdot \text{random.random()} \leftarrow r \in [-4, 0]$   
 $\alpha = 10^r \leftarrow 10^{-4}, \dots, 10^0$

$b = \log_{10} 1$   
 $b = 0 \Rightarrow r \in [a, b] = [-4, 0]$

HPs for exponentially weighted averages

$\beta = 0.9 \dots 0.999$   
 $\beta = 0.9 \rightarrow 10$   
 $\beta = 0.999 \rightarrow 1000$

$\beta = 0.9 \dots 0.999$

→ better not sample uniform randomly  
 → but with log scale!

$1-\beta = 0.1 \dots 0.001$

$r \in [-3, -1]$   
 $1-\beta = 10^r$   
 $\beta = 1-10^r$

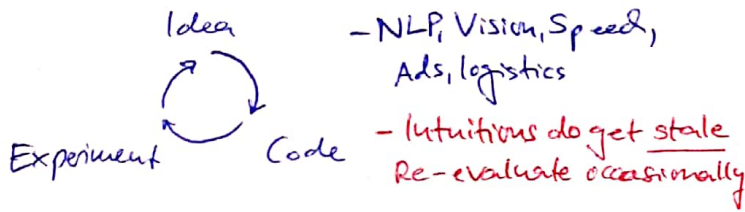
Why is sampling on linear scale a bad idea?

→ when  $\beta$  close to 1 the sensitivity of the results changes even for very low changes in  $\beta$ .

$\beta : 0.9 \rightarrow 0.9005 \sim 10$  values  
 $\beta : 0.999 \rightarrow 0.9995 \sim 1000$  values  
 very sensitive to changes to  $1-\beta$

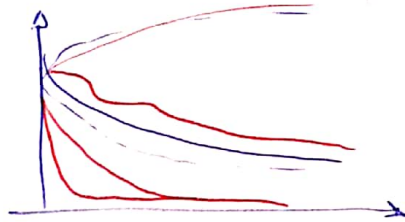
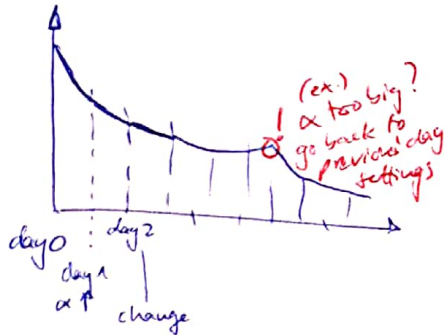
# Pandas vs Caviar

Re-test HPs occasionally



Babysitting one model  
(one model at a time)

Training many models in parallel



"Panda"  
one child at a time

depends also on computational power!

"Caviar"  
many children at a time

Batch normalization → makes HP search easier and the ANN more robust

Normalizing activations in a network → applies not only to inputs but also to deeper layers!

→ Normalizing inputs to speed up learning

$$x_1, x_2, x_3 \xrightarrow{w, b} \hat{y}$$

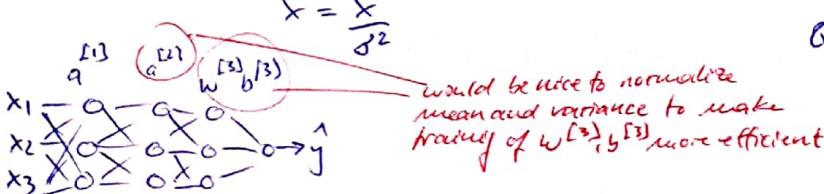
$$\mu = \frac{1}{n} \sum_i x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{n} \sum_i x^{(i)2}$$

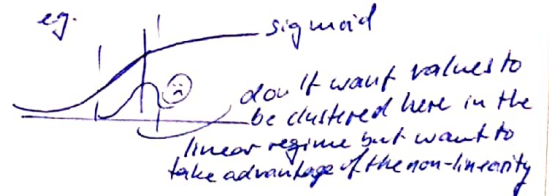
$$x = \frac{x}{\sigma}$$

elementwise



might not want hidden units to be forced to mean 0 and variance 1

eg.



Q: Can we normalize  $a^{(2)}$  so as to train  $w^{(3)}, b^{(3)}$  faster?

→ Normalize  $z^{(2)}$  (before activation but also possible to normalize after activation)

Implement Batch Norm

Given some intermediate values in NN

$$\mu = \frac{1}{n} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{n} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

to have mean 0 and standard unit variance

add  $\epsilon$  for numerical stability if  $\sigma^2 = 0$  for some estimate

$$z^{(1)}, \dots, z^{(n)}$$

$$z^{(i)}$$

but sometimes it might be better for some hidden layers to have a different distribution:

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

learnable params of our model

Update in GD just as  $w, b$

if  $\gamma = \sqrt{\sigma^2 + \epsilon}$  would invert  $\beta = \mu$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

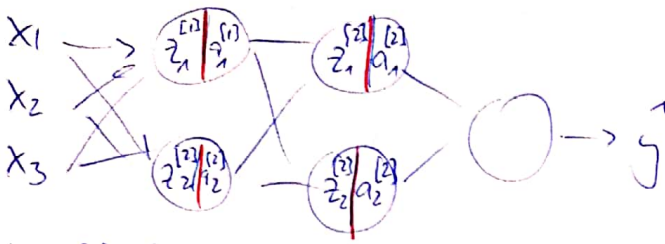
and  $\tilde{z}^{(i)} = z^{(i)}$

Give  $\tilde{z}^{(i)}$  instead of  $z^{(i)}$  for later calculation

②



# Fitting Batch Norm into a NN



In tensor flow: use BN via  
tf.nn.batch-normalization

usually  $x \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \rightarrow a^{[1]}$

with batch norm:

$$x \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \rightarrow a^{[2]}$$

→ mean and variance normalized

Parameters in NN:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$  } Use GD

In practice BN is usually applied with mini-batches

Working with mini-batches

$$x^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \rightarrow \dots$$

then

$$x^{[2]} \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]}$$

$$x^{[3]} \dots x^{[L]}$$

Parameter:  $w^{[L]}, b^{[L]}, \beta^{[L]}, \gamma^{[L]}$

$$\tilde{z}^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

normalize  $\mu$  &  $\sigma$

- because BN zeros out the mean of  $z^{[L]}$   
 $b^{[L]}$  does not need to be kept and is sort of replaced by  $\beta^{[L]}$  → affects shift/bias
- dimension of  $z^{[L]} \rightarrow (n^{[L]}, 1)$   
 $\beta^{[L]}, \gamma^{[L]} \rightarrow (n^{[L]}, 1)$

parameter can be eliminated when using batch norm

compute mean and subtract out the mean so adding a constant is the unbiased to all

examples in a mini-batch doesn't change anything because added constant gets cancelled out

use to decide what is the mean of  $\tilde{z}^{[L]}$

# Implementing gradient descent

assumption: mini-batch is used

for  $t=1 \dots \text{num mini-batches}$

- compute FP on  $X^{t+3}$

In each layer use BN to replace  $z^{(l)}$  with  $\tilde{z}^{(l)}$

→ ensures the mean and variance normalized

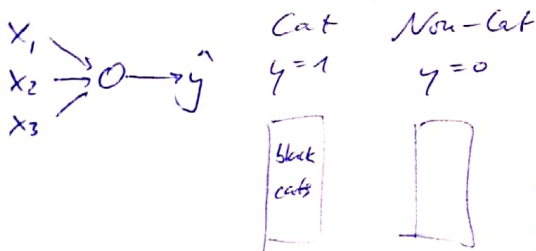
- Use backprop to compute  $dW^{(l)}, db^{(l)}, d\beta^{(l)}, d\gamma^{(l)}$   $z$  is used

- Update parameters  $W^{(l)} := W^{(l)} - \alpha dW^{(l)}$   
 $\beta^{(l)}$   
 $\gamma^{(l)}$

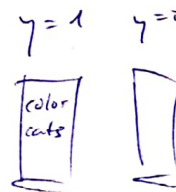
→ works also with GD momentum, RMS prop, Adam

## Why does BN work?

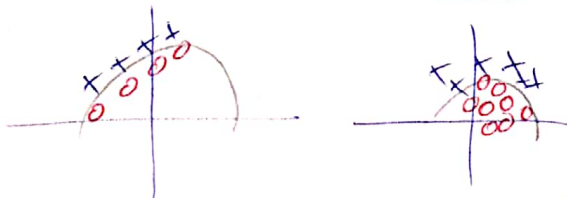
→ makes "later" hidden units more robust to changes in "earlier" hidden units  
 → weakens coupling between layers and so allows each layer to learn by itself more independently of other layers



trained only of images of black cats, now want to apply on images with colored cats  
 classifier will fail



→ speeding up learning!



Would not expect model trained on left data perform well on right data

data distribution changing  
 → Covariate shift

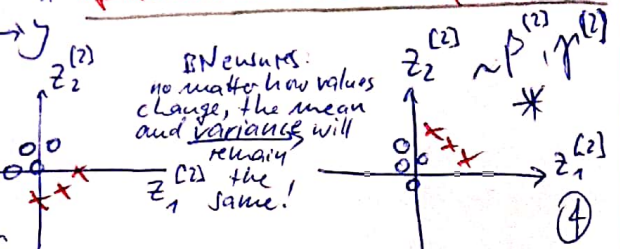
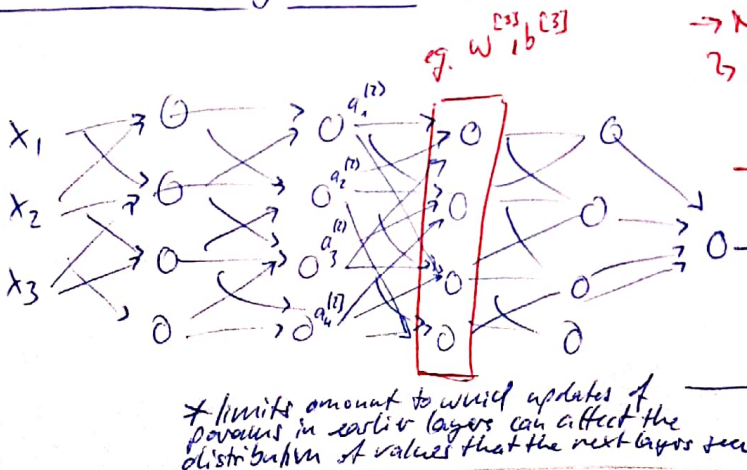
mapping if distribution of  $x$  changes then we might need to re-train our model (ground truth function shift)

How does this apply to a NN? → reduces problem of input values changing

→ NN adapts all params

→ When  $W^{(l)}, b^{(l)}$  change, also  $W^{(l+1)}, b^{(l+1)}$  change  
 → suffering from covariate shift!

→ BN reduces the amount that distribution of hidden unit values shifts around!





Batch Normalization regularization

- Each mini-batch is scaled by mean/variance computed on just that mini-batch
- This adds noise to the values of  $z^{(l)}$  within that mini-batch.
- Similar to dropout, it adds some noise to each hidden layer's activations
- This has slight regularization effect
  - ↳ but not the idea (BN is no regularization method)
- If using a bigger mini-batch size noise is reduced but also regularization effect is reduced

Batch Norm at test time

- BN processes data on mini-batch at a time
- but at test time we might need to process the examples one at a time
- ↳ adapt NN to do that

$$\mu = \frac{1}{m} \sum_i z^{(i)} \quad \text{examples per mini-batch, not whole training set!}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \quad \text{computed on mini-batch}$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

need new way to get  $\mu, \sigma^2$

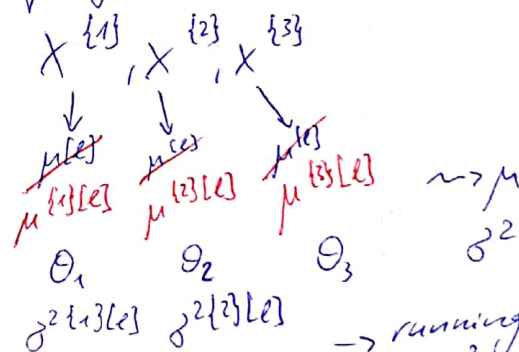
at test time calculate  $z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}$

instead of  $\tilde{z} = \gamma z_{\text{norm}} + \beta$  use exp. weight. avg.

~~find  $\mu, \sigma^2$  on each mini-batch~~

$\mu, \sigma^2$ : estimate via exponentially weighted average (across mini-batches)

e.g. layer  $l = \dots$



→ running average of  $\mu, \sigma^2$  of each layer as we train NN across multiple mini batches

→ estimate  $\mu, \sigma^2$  from training set

→ in practice implement exp. weight. avg. to keep track of  $\mu, \sigma^2$  (via running average)

# Multiclass-classification

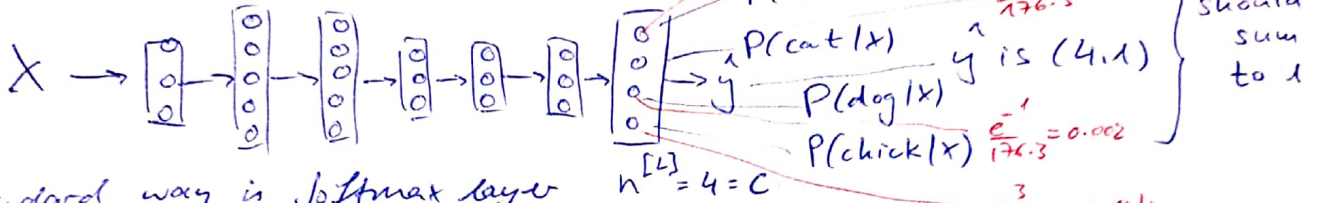
## Softmax Regression

recognizing cats, dogs and baby chicks, others

chick cat dog kowala chick dog kowala cat  
3 1 2 0 3 2 0 1

$C = \# \text{ classes} = 4$

$(0, \dots, 3)$



$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Activation function:

temporary variable  $t = e^{z^{[L]}}$  (4,1)

(element-wise exponentiation)

$$y = a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i}, \quad a_i = \frac{t_i}{\sum_{i=1}^4 t_i}$$

4 classes

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$a^{[L]} = \frac{t}{176.3}$$

$\sum_{i=1}^4 t_i = 176.3$

$$y = a^{[L]} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

$\sum = 1$

$$a^{[L]} = g^{[L]}(z^{[L]})$$

this activation function inputs a (4,1) vector and outputs a (4,1) vector

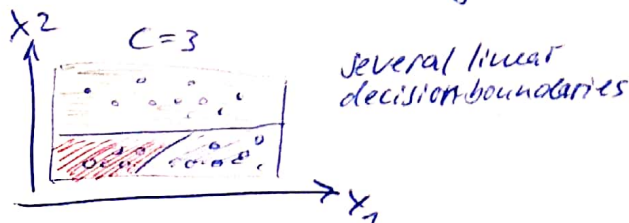
## Softmax examples

no hidden layers

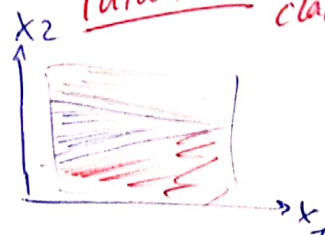
$$x_1 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow y$$

$$z^{[1]} = W^{[1]} x + b^{[1]}$$

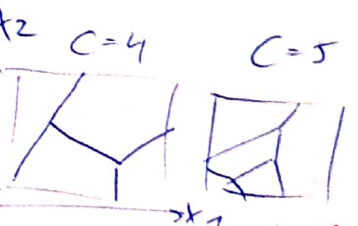
$$a^{[1]} = y = g(z^{[1]})$$



with hidden layers  
also non-linear decision boundaries!!!



Intuition: Decision boundaries will be linear! when there is no hidden layer!



boundaries between classes

# Training a Softmax classifier

example  $z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \rightarrow a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / \dots \\ e^{-1} / \dots \\ e^3 / \dots \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$

Softmax	"hardmax"
$\begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
↓ more gentle mapping	↓ hard mapping

If  $C=2$ , softmax reduces to logistic regression

ex.  $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$  (actually redundant! just compute one and get the other one automatically as they have to add up to 1)

Softmax regression generalizes logistic regression to  $C$  classes

## Loss function

$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  - cat  
 $a^{[L]} = y = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$  - don't look too well

$\mathcal{L}(y^i, y) = - \sum_{j=1}^4 y_j \log(y_j^i)$

$= -y_2 \log y_2$  (all except "2" were 0)

to make  $\mathcal{L}(y^i, y)$  small we need to make  $-\log y_2$  small, so we make  $y_2$  big!

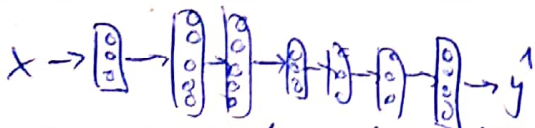
$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^i, y^{(i)})$

use GD to minimize  $J$

$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$      $Y^a = [y^{a(1)}, \dots, y^{a(m)}]$

$= \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}_{(4,m)}$      $= \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}_{(4,m)}$

## Gradient descent with softmax



FP is straight forward  $z^{[L]} \rightarrow a^{[L]} = g \rightarrow \mathcal{L}(y, y)$

BP:  $\frac{\partial \mathcal{L}}{\partial z^{[L]}} = g' - y$



# Deep Learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

## Choosing DL Framework

- Ease of programming (development and deployment)
- Running speed
- Truly open (open source with good governance)

## Tensor Flow (1) / Programming assignment

ex.  $J(w) = |w^2 - 10w + 25|$

create & use sessions  
method 1

```
sess = tf.Session()
result = sess.run(..., feed_dict={w: 5})
sess.close()
```

method 2

with `tf.Session()` as sess:

feed data: `result = sess.run(tf.sigmoid(x), feed_dict={x: 2})`  
for running a session (closes session automatically)

dictionary here

with arrays as inputs  $\rightarrow$  `result = sess.run(cost, feed_dict={x: logits, y: labels})`

"one hot" encoding  
result vector

$y = [1, 2, 0, 3, 2]$   $\rightarrow$   $C=3$

0	0	1	0	0	class 0
1	0	0	0	0	class 1
0	1	0	0	1	class 2
0	0	0	1	0	class 3

$\rightarrow$  `tf.one_hot(labels, depth=C, axis=-1)`

$\rightarrow$  `tf.one_hot(indices=labels, depth=C, axis=-1)`

initialize to zero and one

`tf.ones(shape)`  
`tf.zeros(shape)`

$\rightarrow$  Create tensors/variables that are not executed yet

$\rightarrow$  Define operations between the tensors  
*tf.matmul / tf.add*

$\rightarrow$  Initialize tensors

$\rightarrow$  Create a session

$\rightarrow$  run the session

TensorFlow constant with value 5

`tf.constant(5, name="y")`

`a = tf.placeholder(tf.float32, name="a")`

`x = tf.variable(a + 3, name="x")`

constant or placeholder



# Tensor Flow / Programming assignment sign language!

- Create a computation graph
- Run the graph

train set = 1080 pictures (64x64 pixels) RGB  
test set = 120 pictures (64x64x3)  
C = 6 → hand signals from 0 to 5

- flatten the image + normalize by 255

$X_{train\_flatten} = X_{train\_orig}.reshape(X_{train\_orig}.shape[0], -1).T$

↳ same for Y and test...

- Create placeholders

example  $x = tf.placeholder(tf.float32, shape = (...), name = "x")$  or  $[..., ...]$

- initialize parameters

example  $w1 = tf.get_variable("w1", [25, 12288], initializer = tf.contrib.layers.xavier_initializer)$

- Forward propagation

→ outputs no cache here...

example  $z2 = tf.add(tf.matmul(w2, A1), b2)$   
 $A2 = tf.nn.relu(z2)$

- Compute cost

input(z3, Y) → logits = tf.transpose(z3)  
labels = tf.transpose(Y)

$cost = tf.reduce\_mean(tf.nn.softmax\_cross\_entropy\_with\_logits(logits = logits, labels = labels))$

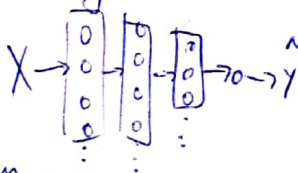
- Backpropagation + parameter updates

→ after cost, create object optimizer, called together with cost in tf.session

eg. optimizer = tf.train.GradientDescentOptimizer(learning\_rate = learning\_rate).minimize(cost)  
to do optimization:  $optimizer = tf.train.AdamOptimizer(...)$

→  $c = sess.run([optimizer, cost], feed_dict = \{X: mini\_batch\_X, Y: mini\_batch\_Y\})$

- Building the model



$$X = n_x \begin{bmatrix} | & | & | & | \\ | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{bmatrix}$$