

TensorFlow

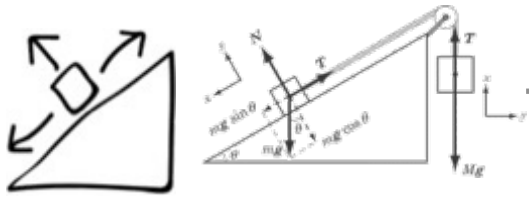
symbolic computing for machine learning



this is about creating models



ANALYTIC MODELS

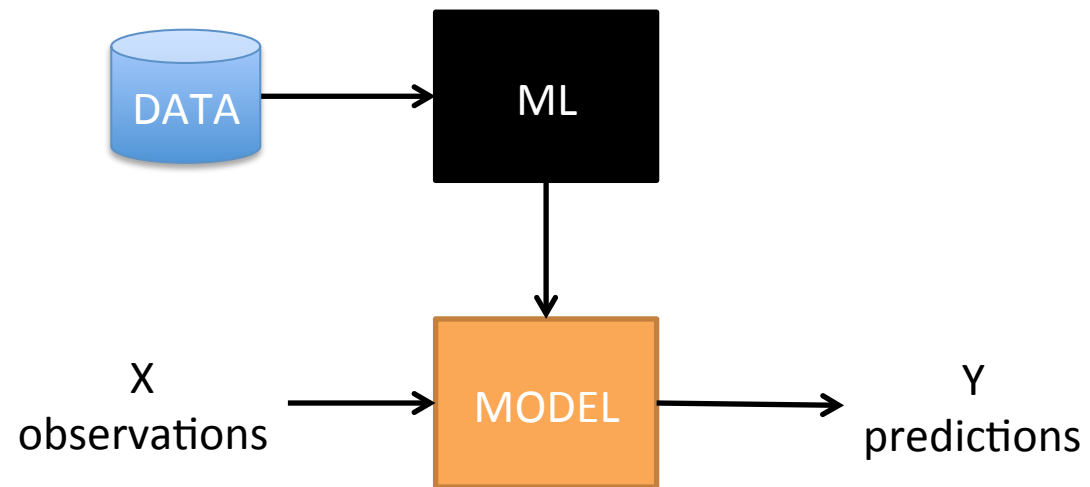


X
observations

$$\begin{aligned}\frac{1}{2}Mv^2 + \frac{1}{2}mv^2 - Mgh + mgh \sin \theta &= 0 \\ \frac{1}{2}(m + M)v^2 &= gh(M - m \sin \theta) \\ v &= \sqrt{\frac{2gh(M - m \sin \theta)}{m + M}}\end{aligned}$$

Y
predictions

MACHINE LEARNING



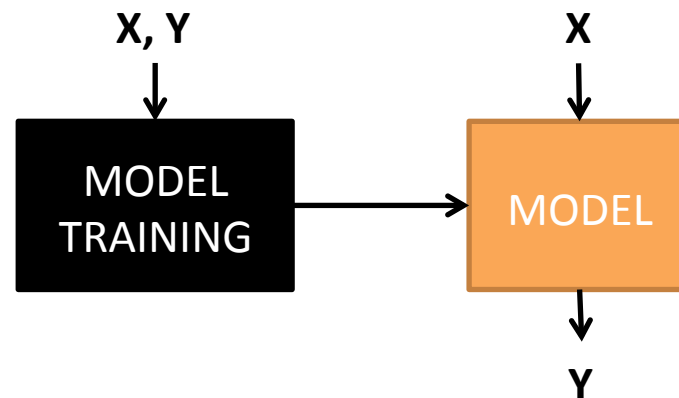
SUPERVISED LEARNING

X

Patient data
Images
Financial data
Robot sensors
Flight history
Movies watched
Recorded voice
Driving data

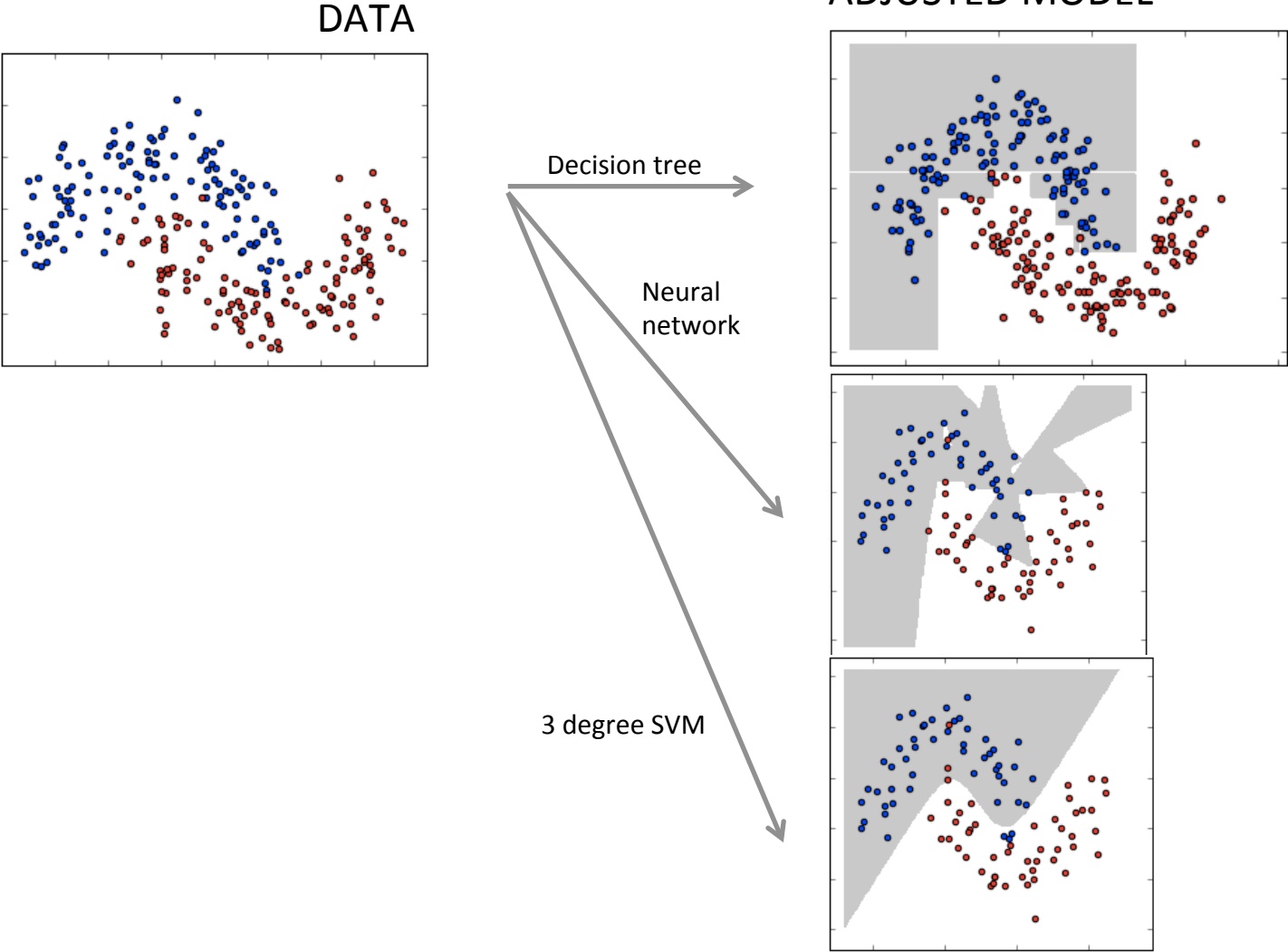
Y

Presence of pathology
Objects present
Next price
Steering action
Landing time
Movies to suggest
Words transcribed
Driver profile

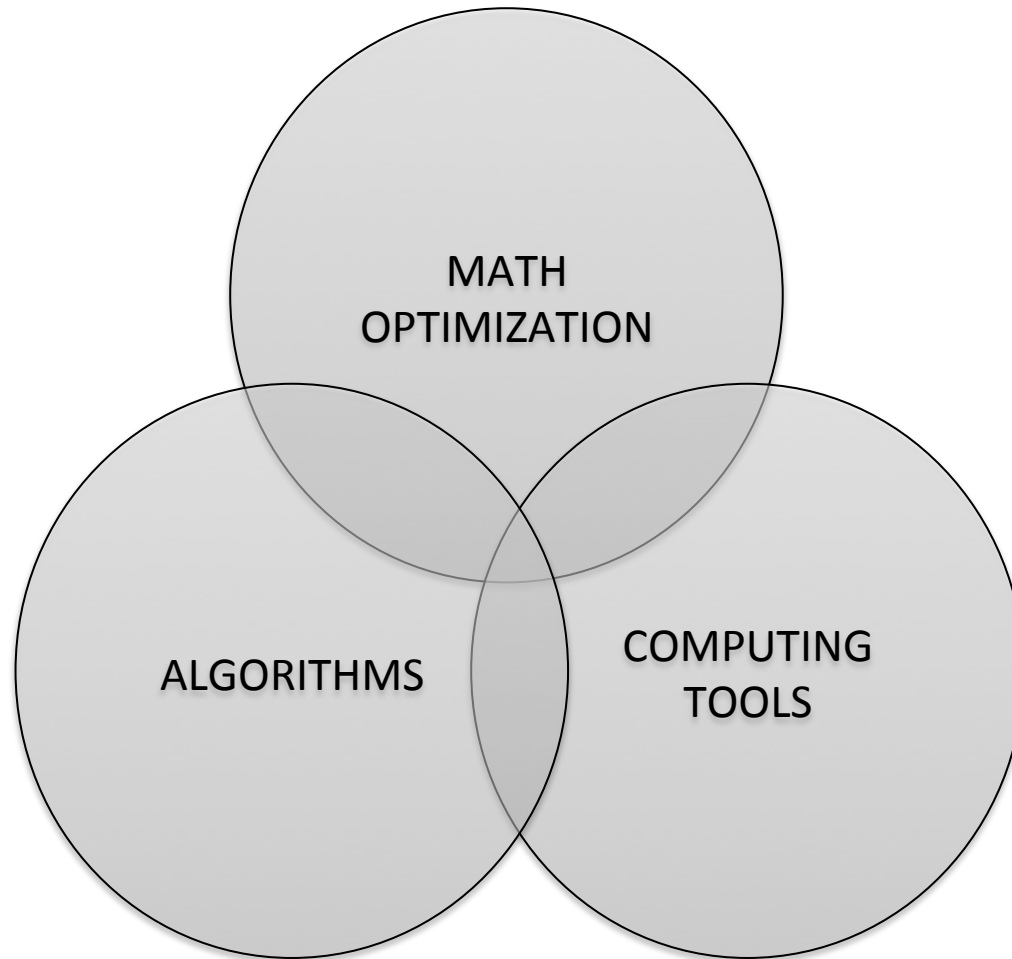


show video

FAMILIES of MODELS



MACHINE LEARNING



STANDARD FORMULATION

Tenemos un dataset de entrenamiento $\{(x^{(i)}, y^{(i)})\}$ con $i \in \{1 \dots m\}$ (es decir, con m puntos), $x^{(i)} = [1 \ x_1^{(i)} \ x_2^{(i)} \ \dots \ x_j^{(i)} \ \dots \ x_n^{(i)}]^T \in \mathbb{R}^{n+1}$ y $y^{(i)} \in \{0, 1\}$ para clasificación binaria. La matriz \mathbf{X} recoge todos los $x^{(i)}$ y el vector \mathbf{y} todos los $y^{(i)}$

$$\mathbf{X} = \begin{bmatrix} \text{---} & x^{(1)} & \text{---} \\ \text{---} & x^{(2)} & \text{---} \\ & \dots & \\ \text{---} & x^{(m)} & \text{---} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ & & \dots & & \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\hat{y}^{(i)} = f(x^{(i)}) \quad \text{where } f \text{ is our model}$$

m : número de datos n : número de descriptores por dato

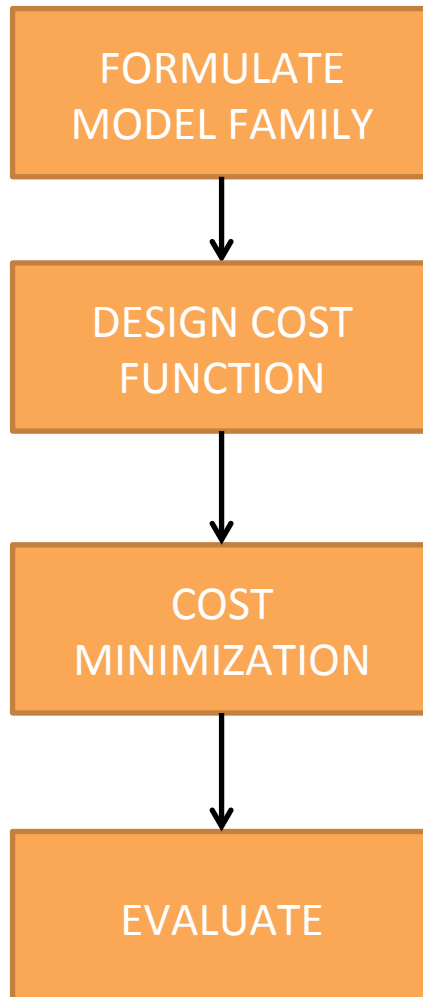
STANDARD FORMULATION

Tenemos un dataset de entrenamiento $\{(x^{(i)}, y^{(i)})\}$ con $i \in \{1 \dots m\}$ (es decir, con m puntos), $x^{(i)} = [1 \ x_1^{(i)} \ x_2^{(i)} \ \dots \ x_j^{(i)} \ \dots \ x_n^{(i)}]^T \in \mathbb{R}^{n+1}$ y $y^{(i)} \in \{0, 1\}$ para clasificación binaria. La matriz \mathbf{X} recoge todos los $x^{(i)}$ y el vector \mathbf{y} todos los $y^{(i)}$

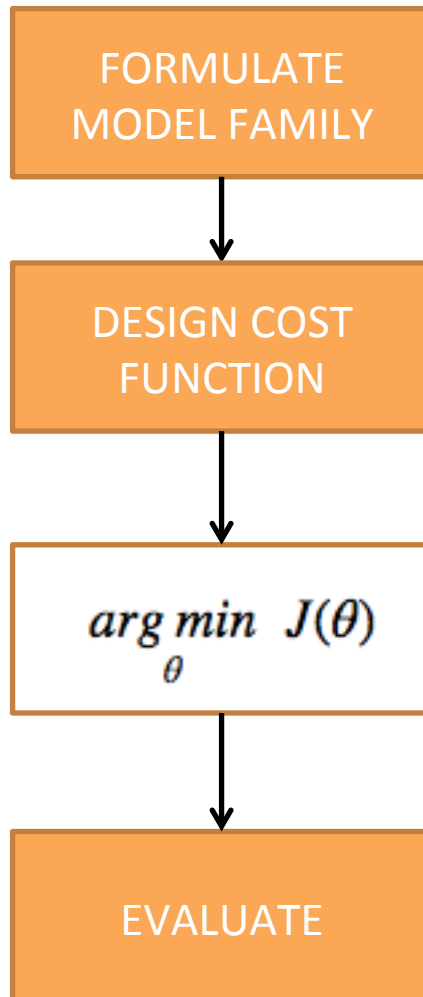
$$\mathbf{X} = \begin{bmatrix} \text{---} & x^{(1)} & \text{---} \\ \text{---} & x^{(2)} & \text{---} \\ & \dots & \\ \text{---} & x^{(m)} & \text{---} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ & & \dots & & \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$y^{(i)} \in \mathbb{R} \quad \rightarrow \quad \textit{Regression}$$
$$y^{(i)} \in [0, 1] \textit{ or } [-1, 1] \quad \rightarrow \quad \textit{Classification}$$

STANDARD WORKFLOW



LINEAR REGRESSION



X/Y linear relation

$$\hat{y}^{(i)} = \theta^T x^{(i)}$$

Avg squared error

$$J(\theta) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2$$

Gradient descent
(or generic optimizer)

$$\nabla J = \frac{2}{m} X^T \cdot (X \cdot \theta - Y)$$

Closed form

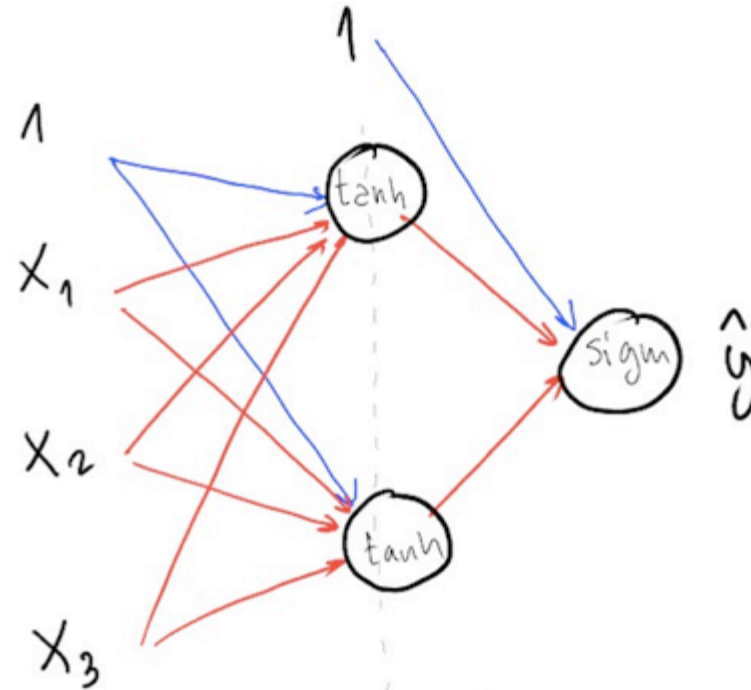
$$\theta = (X^T X)^{-1} X^T Y$$

Train/test split
Learning curves,
Stability of models, etc.

show notebook part 1

PERCEPTRON

FORMULATE
MODEL FAMILY



$$\mathbf{W}_1 = \begin{bmatrix} w_0^0 & w_1^0 \\ w_0^1 & w_1^1 \\ w_0^2 & w_1^2 \end{bmatrix}$$

$$\mathbf{b}_1 = [b_0 \quad b_1]$$

params: b_1 W_1 b_2 W_2

dims: 1×2 3×2 1 2×1

$$\hat{y} = \text{sigm} (W_2 \cdot \tanh (x^T W_1 + b_1) + b_2)$$

PERCEPTRON

FORMULATE
MODEL FAMILY



DESIGN COST
FUNCTION



$\arg \min_{\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2} J(\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2)$

$$\hat{y} = g(\mathbf{W}_2 \cdot \tanh(\mathbf{x}^T \cdot \mathbf{W}_1 + \mathbf{b}_1) + b_2)$$

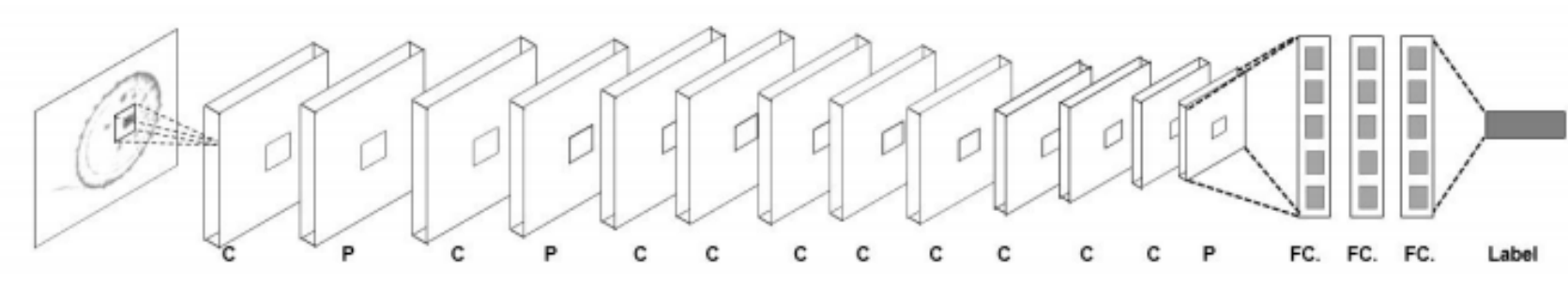
$$J(\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2 .$$

Gradient descent (or generic optimizer)

BACKPROPAGATION

CONVOLUTIONAL NETWORKS

ml4a.github.io/dev/demos/demo_convolution.html



STRUCTURE OF COST FUNCTIONS

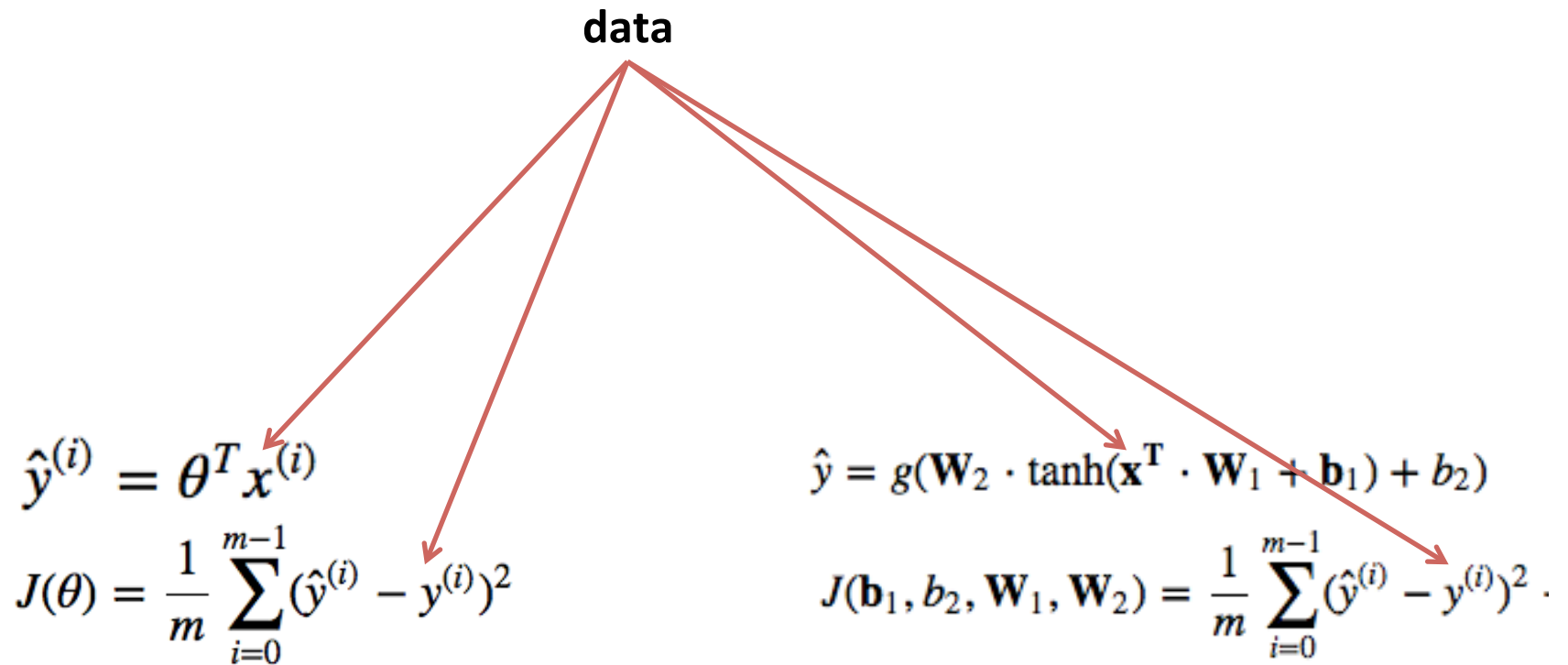
$$\hat{y}^{(i)} = \theta^T x^{(i)}$$

$$J(\theta) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2$$

$$\hat{y} = g(\mathbf{W}_2 \cdot \tanh(\mathbf{x}^T \cdot \mathbf{W}_1 + \mathbf{b}_1) + b_2)$$

$$J(\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2 .$$

STRUCTURE OF COST FUNCTIONS



STRUCTURE OF COST FUNCTIONS

params



$$\hat{y}^{(i)} = \theta^T x^{(i)}$$

$$J(\theta) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2$$

$$\hat{y} = g(\mathbf{W}_2 \cdot \tanh(\mathbf{x}^T \cdot \mathbf{W}_1 + \mathbf{b}_1) + b_2)$$

$$J(\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2 .$$

STRUCTURE OF COST FUNCTIONS

params

**MATRIX SYMBOLIC COMPUTING
WITH SYMBOLS FOR DATA (no gradients)
AND SYMBOLS FOR PARAMS (need gradients)**

$$\hat{y}^{(i)} = \theta^T x^{(i)}$$

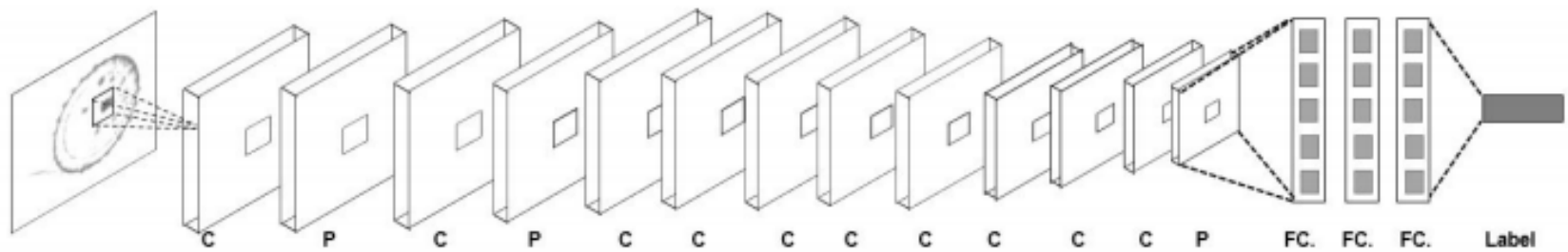
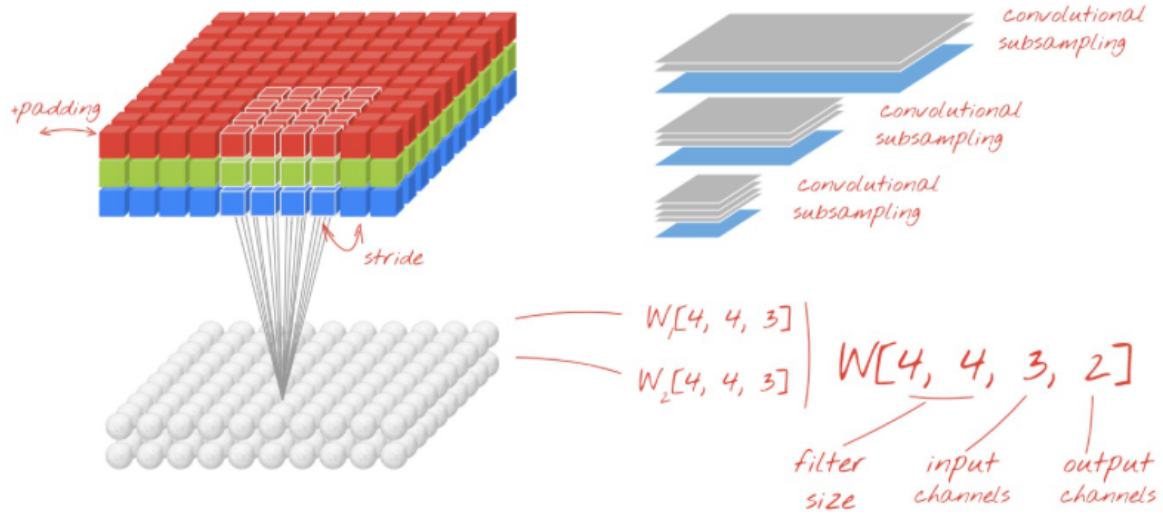
$$J(\theta) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2$$

$$\hat{y} = g(\mathbf{W}_2 \cdot \tanh(\mathbf{x}^T \cdot \mathbf{W}_1 + \mathbf{b}_1) + b_2)$$

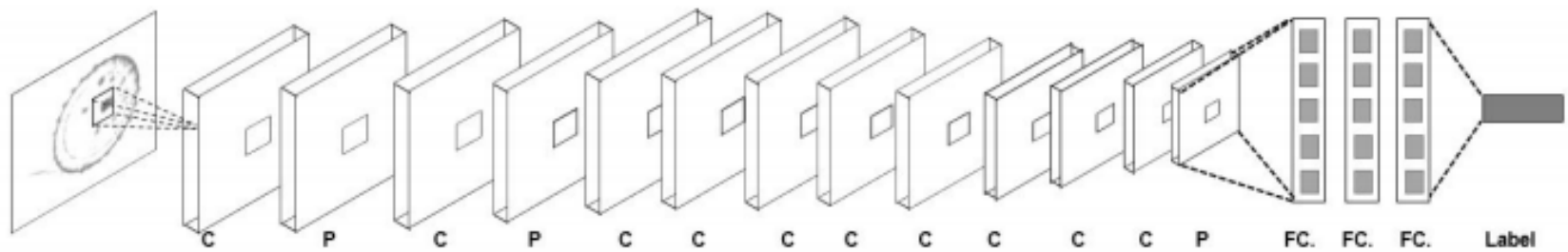
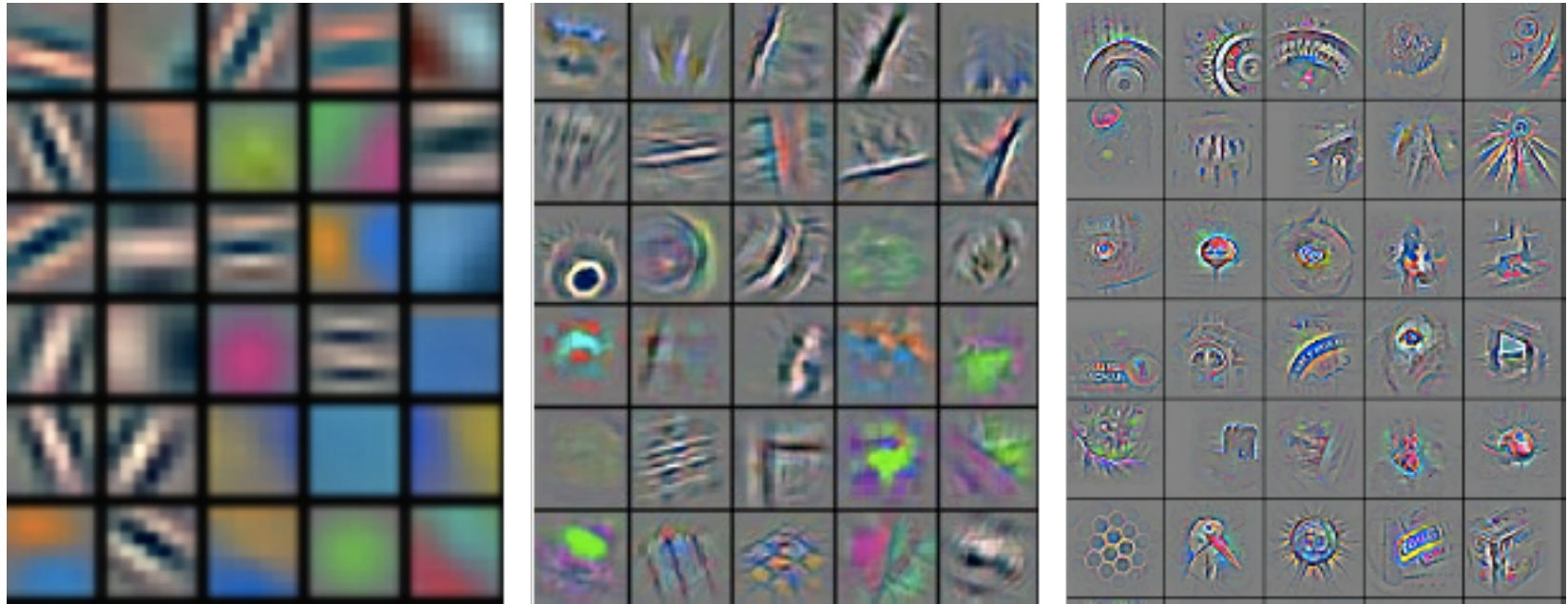
$$J(\mathbf{b}_1, b_2, \mathbf{W}_1, \mathbf{W}_2) = \frac{1}{m} \sum_{i=0}^{m-1} (\hat{y}^{(i)} - y^{(i)})^2 .$$

show notebook part 2

CNN



CNN



“small” CNN



layer	input_size	output_size	filter_size	stride	n_filters	activation	W_size from previous
conv1	28x28x1	28x28x9	5x5	1	16	relu	W1 = [5,5,1,16]
conv2	28x28x16	14x14x8	5x5	2	8	relu	W2 = [5,5,16,8]
conv3	14x14x8	7x7x12	4x4	2	12	relu	W2 = [4,4,8,12]
fc	7x7x12	200				relu	W3 = [588,200]
output	200	10				softmax	W4 = [200,10]

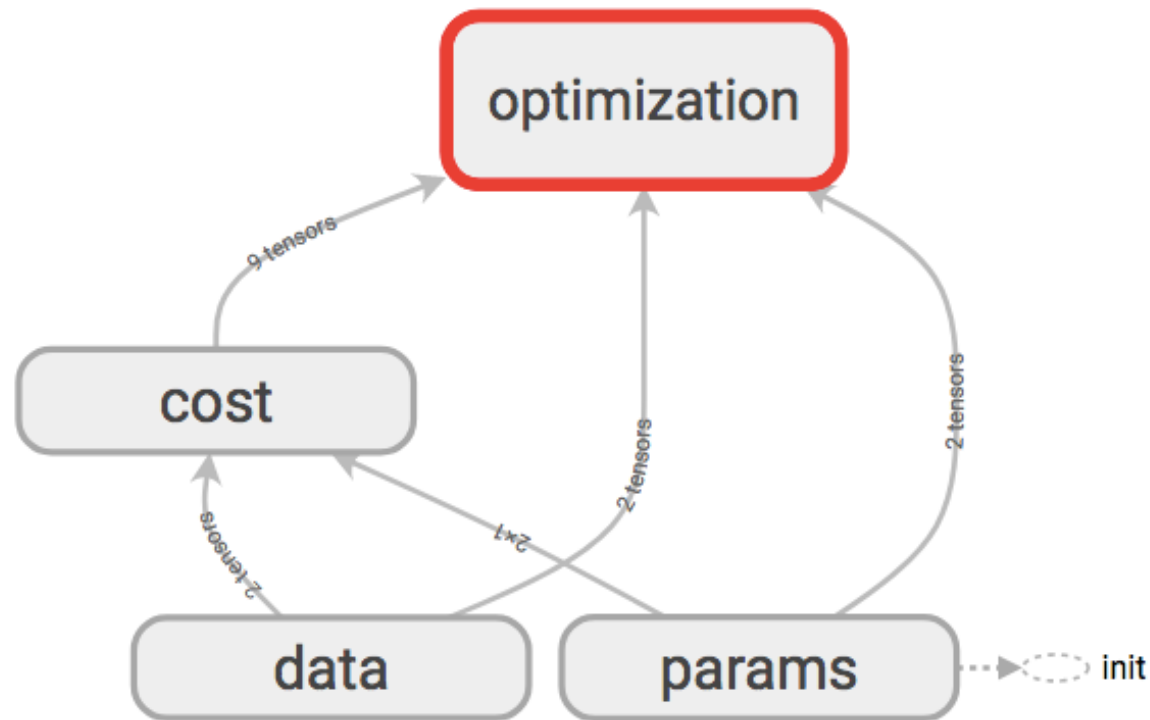
$400 + 3600 + 1536 + 117600 + 2000 = 125136$ tunnable params!!!

show notebook Notes 04

show logdir in tensorboard

device placement

show graph in tensorboard from last CNN



device placement

- individual ops have parallel implementations (multi core CPU or multi thread GPU)
- can specify device placement of components of the computation graph (data and/or operations)
- tensorflow will move data around to comply

```
# Creates a graph.
with tf.device('/device:GPU:2'):
    a = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[2, 3], name='a')
    b = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[3, 2], name='b')
    c = tf.matmul(a, b)
# Creates a session with log_device_placement set to True.
sess = tf.Session(config=tf.ConfigProto(log_device_placement=True))
# Runs the op.
print(sess.run(c))
```

device placement

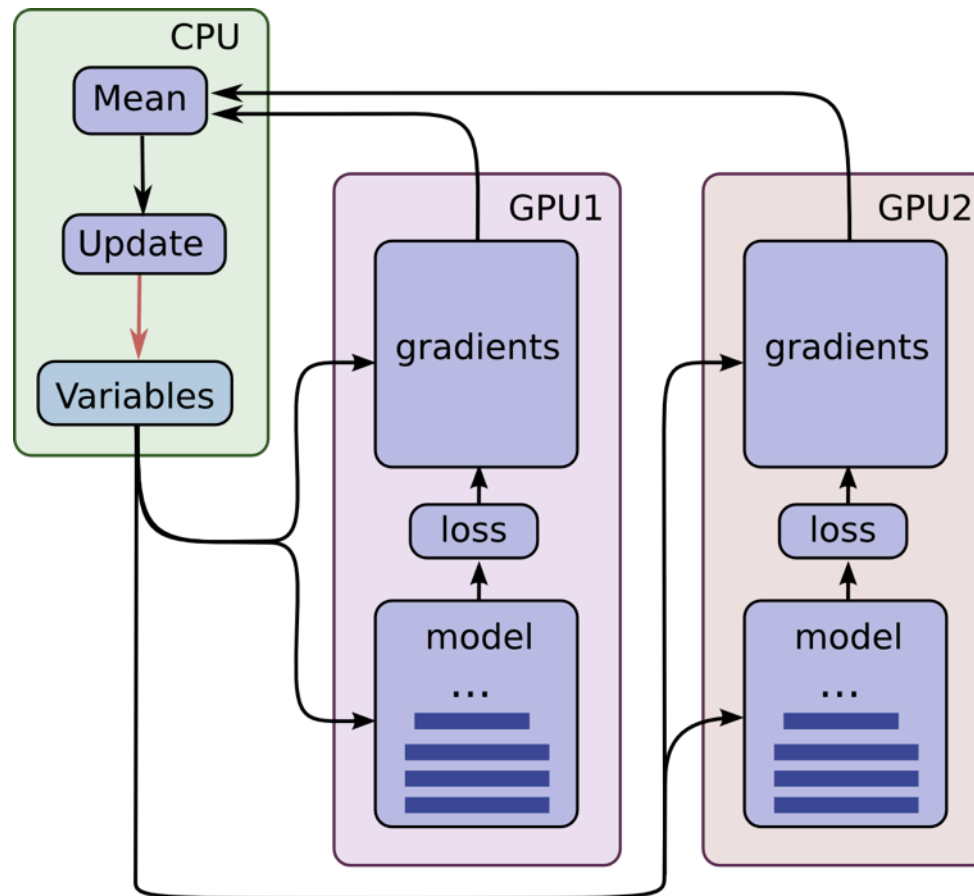
- types of devices
- not all operations can be done in GPU (ops vs. kernels)

```
# Creates a graph.
c = []
for d in ['/device:GPU:2', '/device:GPU:3']:
    with tf.device(d):
        a = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[2, 3])
        b = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[3, 2])
        c.append(tf.matmul(a, b))
with tf.device('/cpu:0'):
    sum = tf.add_n(c)
# Creates a session with log_device_placement set to True.
sess = tf.Session(config=tf.ConfigProto(log_device_placement=True))
# Runs the op.
print(sess.run(sum))
```

recall that GPUs are co-processors, everything exists first in RAM

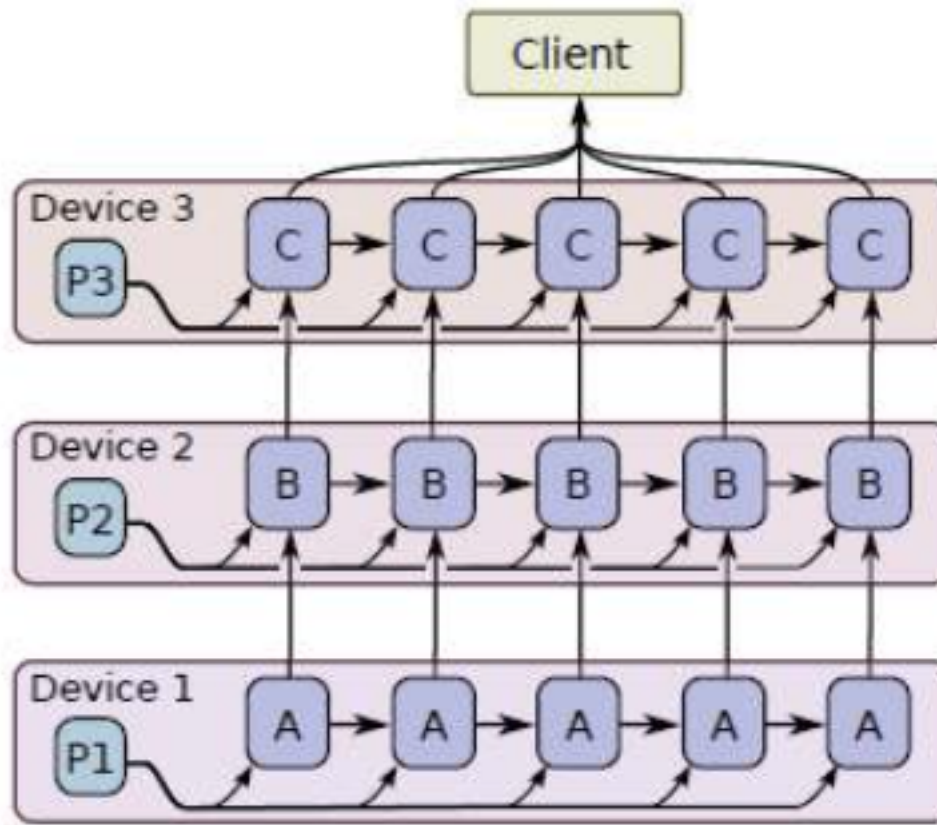
data parallelism

- data parallelism through batches of data
- model + batch needs to fit in GPU memory



model parallelism

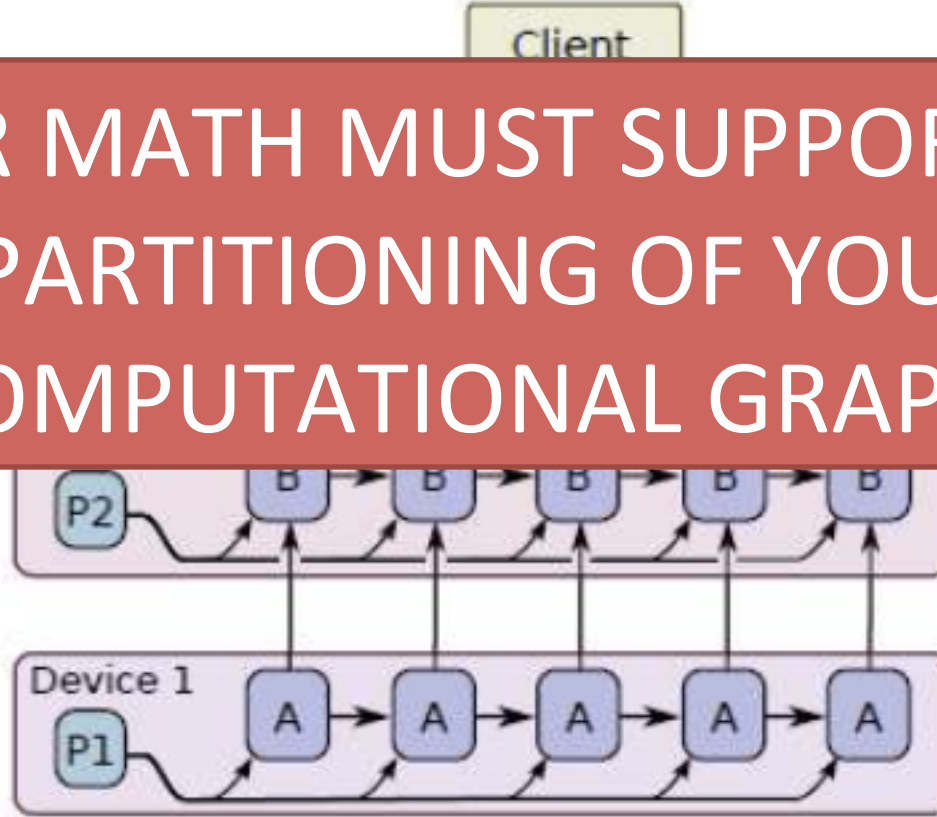
- harder to program
- needs sync
- GPU mem holds partial model and data batches



model parallelism

- harder to program
- needs sync
- GPU mem holds partial model and data batches

YOUR MATH MUST SUPPORT THE
PARTITIONING OF YOUR
COMPUTATIONAL GRAPH!!



distributed computing

data is expensive to move!!!!

in/between graph replication, asynch/synch, etc.

```
with tf.device("/job:ps/task:0"):
    weights_1 = tf.Variable(...)
    biases_1 = tf.Variable(...)

with tf.device("/job:ps/task:1"):
    weights_2 = tf.Variable(...)
    biases_2 = tf.Variable(...)

with tf.device("/job:worker/task:7"):
    input, labels = ...
    layer_1 = tf.nn.relu(tf.matmul(input, weights_1) + biases_1)
    logits = tf.nn.relu(tf.matmul(layer_1, weights_2) + biases_2)
    # ...
    train_op = ...

with tf.Session("grpc://worker7.example.com:2222") as sess:
    for _ in range(10000):
        sess.run(train_op)
```

distributed computing

data is expensive to move!!!!

in/between graph replication, asynch/synch, etc.

```
with tf.device("/job:ps/task:0"):
    weights_1 = tf.Variable( ... )
```

THINK ON HOW DATA FLOWS
WITHIN YOUR COMPUTATIONAL GRAPH

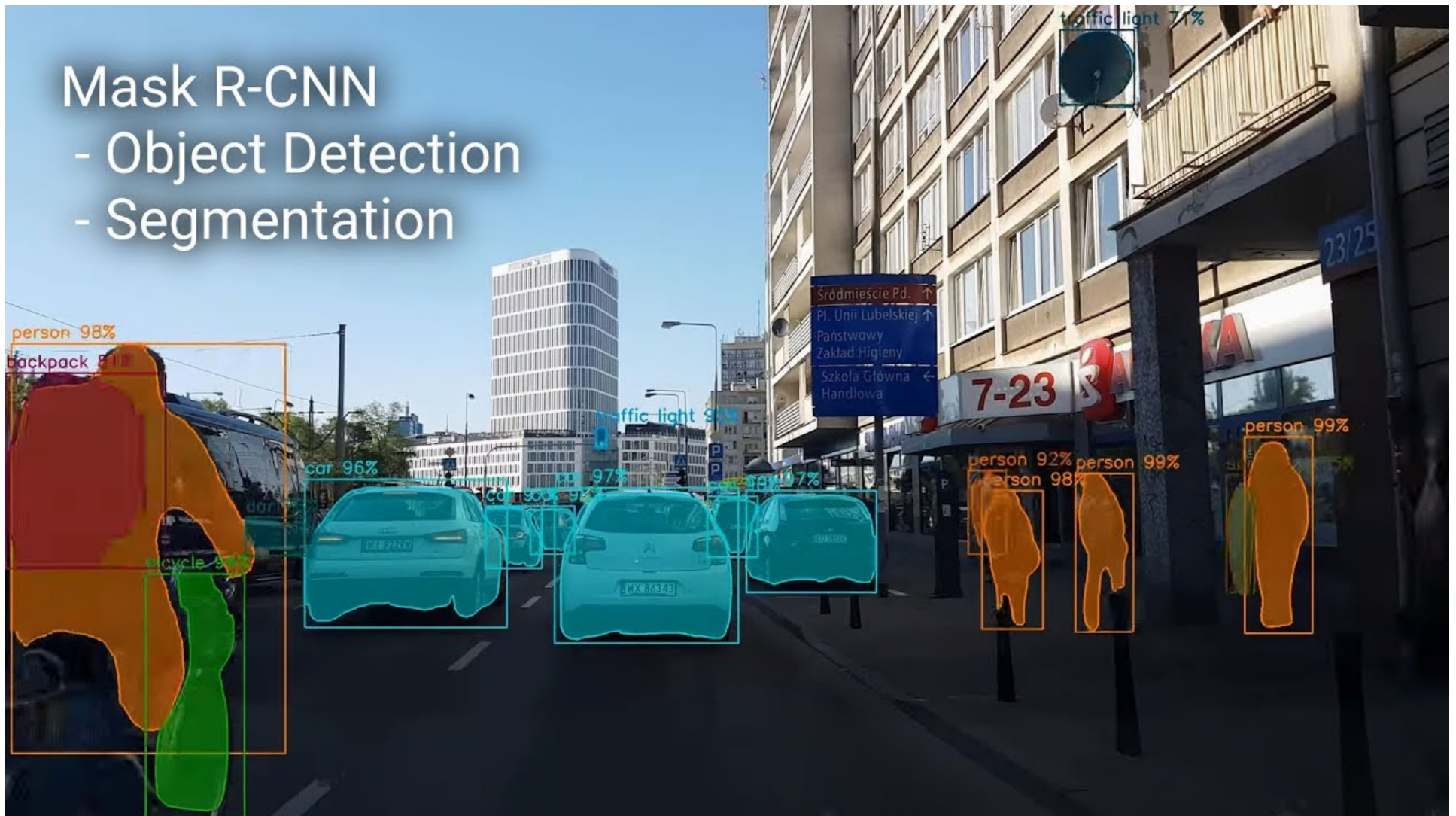
```
with tf.device("/job:worker/task:7"):
    input, labels = ...
    layer_1 = tf.nn.relu(tf.matmul(input, weights_1) + biases_1)
    logits = tf.nn.relu(tf.matmul(layer_1, weights_2) + biases_2)
    # ...
    train_op = ...

with tf.Session("grpc://worker7.example.com:2222") as sess:
    for _ in range(10000):
        sess.run(train_op)
```

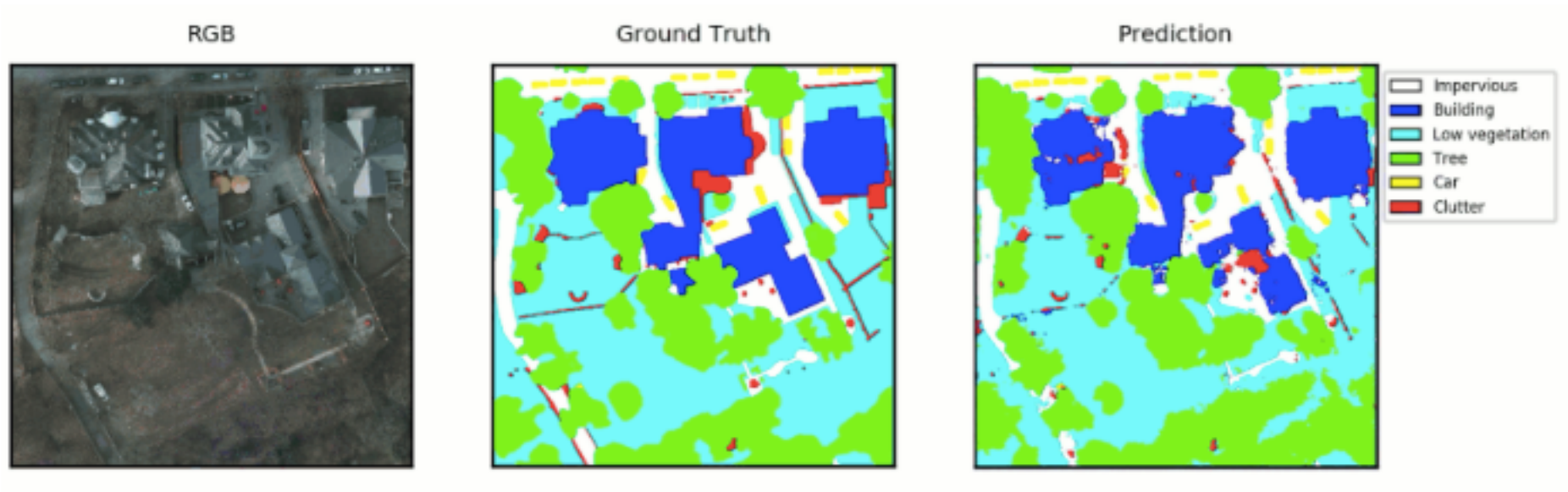

CNN applications

Mask R-CNN

- Object Detection
- Segmentation

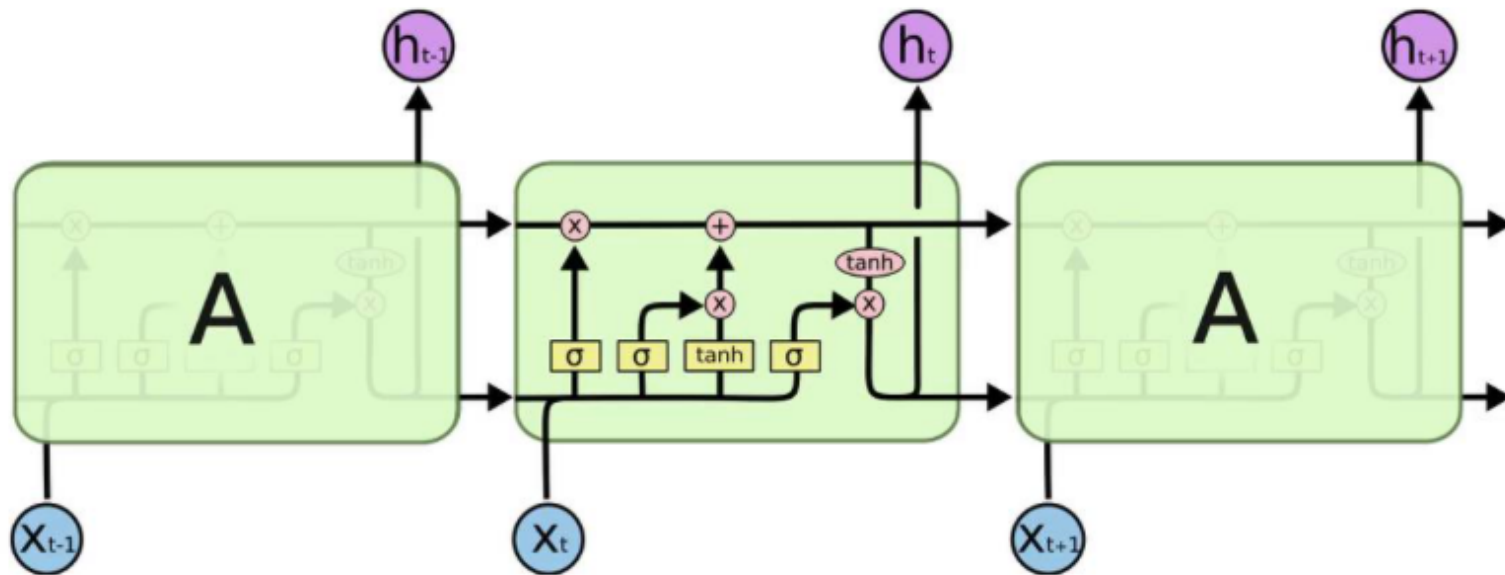


CNN applications

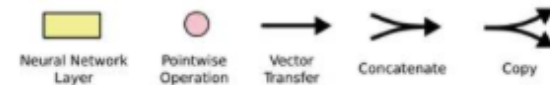


RNNs

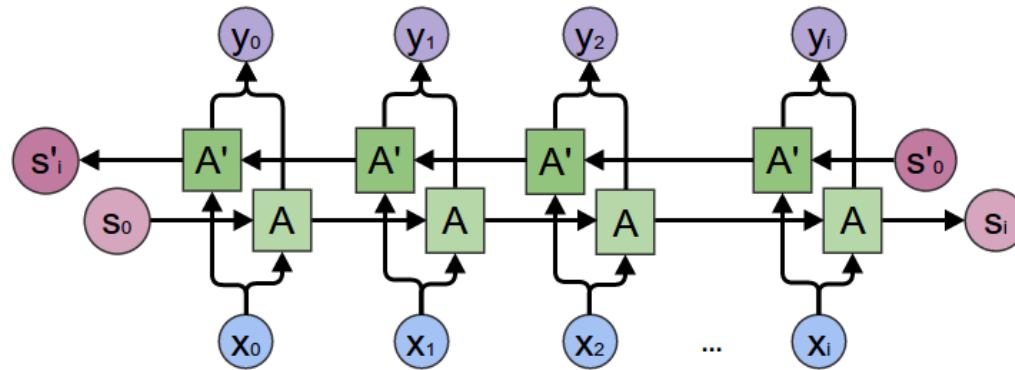
Long-Short Term Memory module: LSTM



long-short term memory modules used in an RNN



RNN applications



signal patterns (finance, speech, etc.)
text generation, translation

see <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

other stuff

- High Level API (tflearn, Keras)
- TPUs (tensor processing units)
- Theano and Torch

https://s3.amazonaws.com/rlx/streetview_detection/madrid_centro/dboard.html

https://s3.amazonaws.com/rlx/eafit_edificios/dboard_clean.html

RISE study group on TF

- Learn and discuss on TF
- Understand its applicability
- Approach problems and practical solutions
- Compete in challenges!!!!

- Open to anybody
- Python programming encouraged

<https://goo.gl/Vf6v4A>

thnx

