



DEPARTAMENTO DE SEÑALES, SISTEMAS Y RADIOCOMUNICACIONES



Deep Learning Seminar Day-1

Master of Science in Signal Theory and Communications
TRACK: Signal Processing and Machine Learning for Big Data

Departamento de Señales, Sistemas y Radiocomunicaciones
E.T.S. Ingenieros de Telecomunicación
Universidad Politécnica de Madrid

Deep Learning Seminar Overview

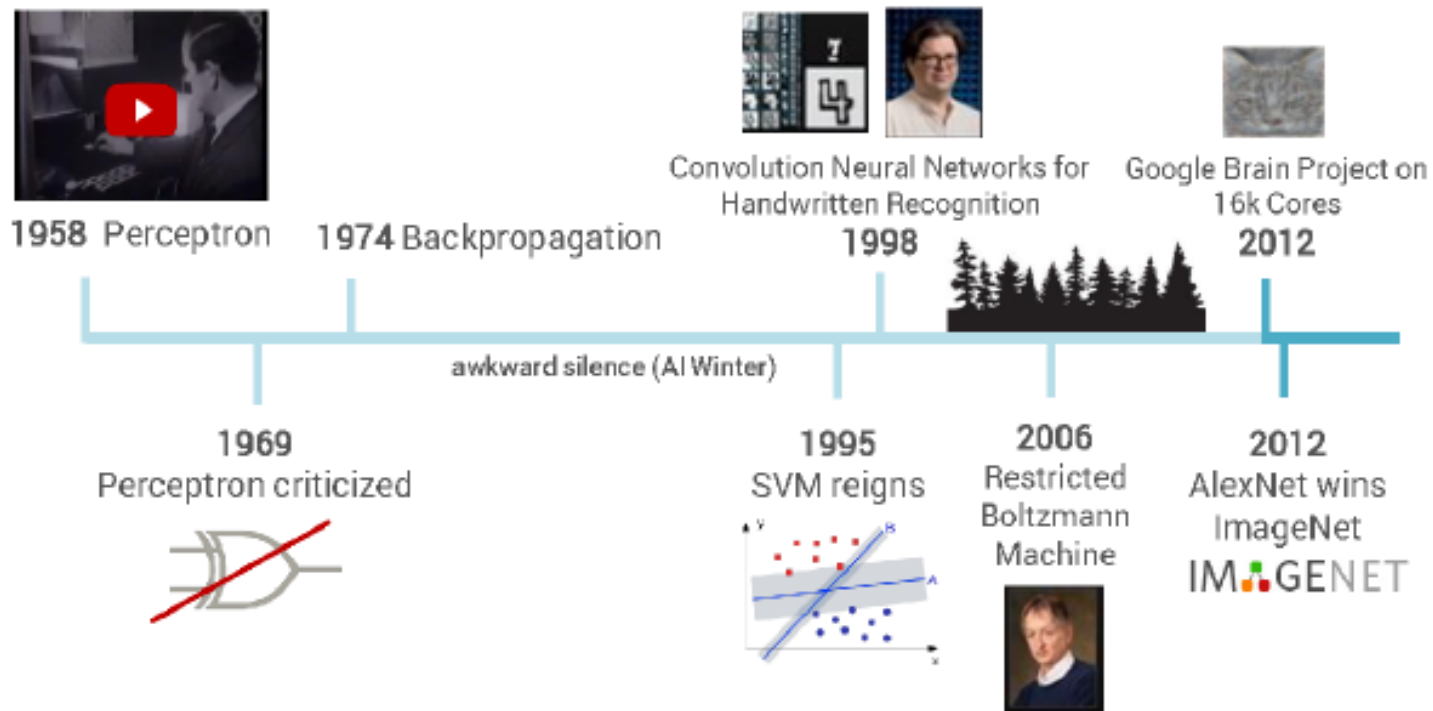
- **Day-1:** Introduction to Deep Learning and Tensorflow
 - Introducing Deep Neural Networks from Linear Classifier (logistic regression)
 - Gradient Descent
 - Simple use of Tensorflow for binary linear classification using artificially generated data
- **Day-2:** Deep Learning Architectures
 - Backpropagation overview and Deep Learning models
 - Using Tensorflow for Image Classification: Logistic Regression, Feed Forward and Convolutional Networks (CNN)
- **Day-3:** Recurrent Neural Networks
 - RNN fundamentals, truncated backpropagation
 - Using Vanilla RNN, LSTM, GRU in Tensorflow
 - Simple Natural Language Processing examples
 - Best-practice discussion

Deep Learning: Hype or Reality?



Deep Learning: Hype or Reality?

A Brief History



Deep Learning using TensorFlow and TensorFlow-Slim

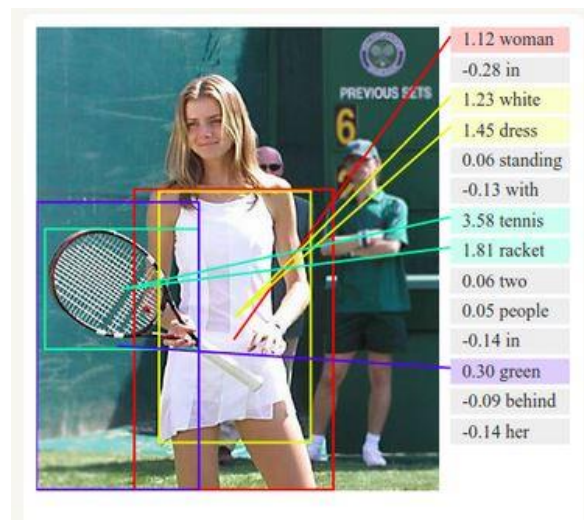
Dipendra Jha Northwestern University

Deep Learning: Hype or Reality?

DNNs better than humans at image recognition

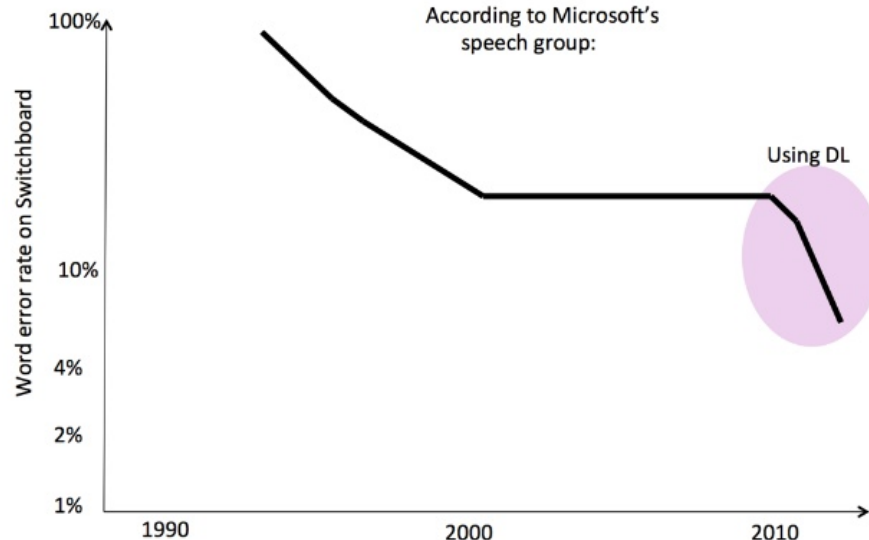


ImageNet
1000 categories
1.3M images
Human error: 5%
DNN: 3%



Speech Recognition

According to Microsoft's speech group:



Deep Learning: Hype or Reality?

Just with in Google

- Search
- Search by image
- Driveless cars
- Youtube recommendation
 - Videos
 - Thumbnails
- Maps
 - Reading street addresses



Deep Learning: Hype or Reality?

Deep Reinforcement Learning



...going unsupervised! Deep Clustering

Introduction to Deep Learning

but what is new?

What Changed?
Old wine in new bottles



Big Data
(Digitalization)



Computation
(Moore's Law, GPUs)

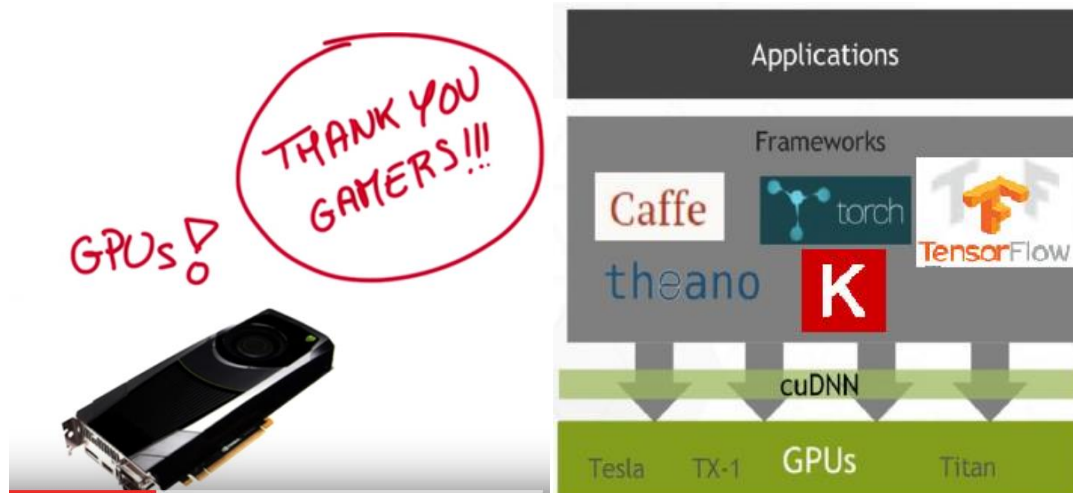


Algorithmic
Progress

UDACITY

FREE COURSE

Deep Learning
by Google

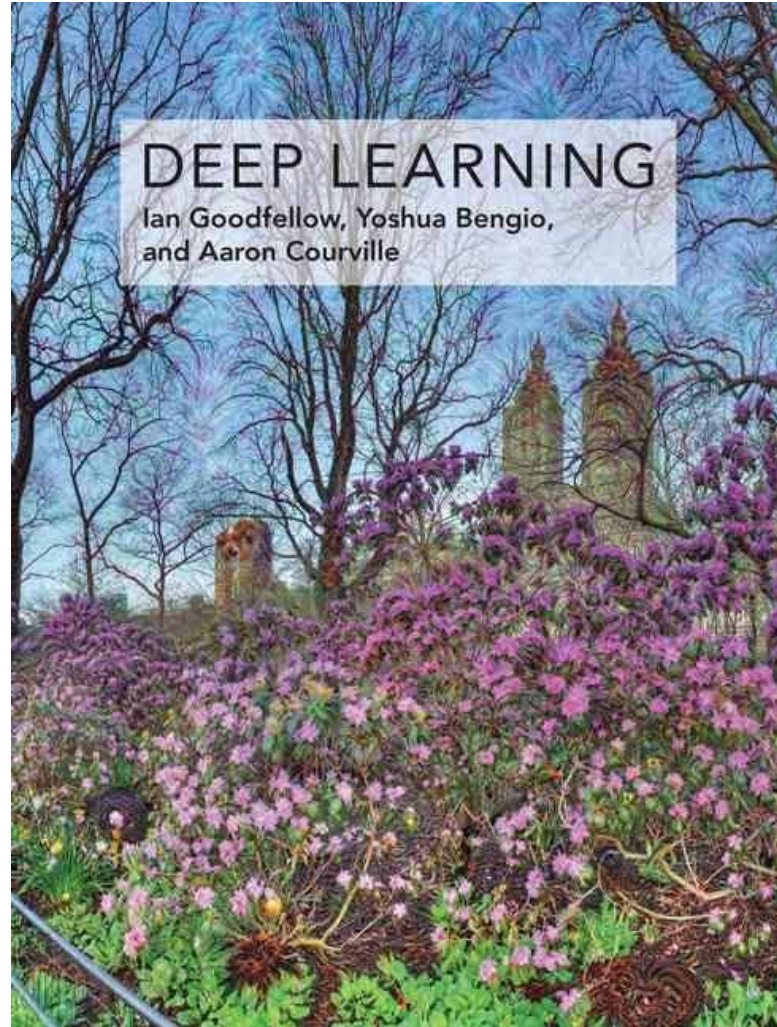


Introduction to Deep Learning

...the Machine Learning background...

Learn the whole
Machine Learning
context

On line:
www.deeplearningbook.org

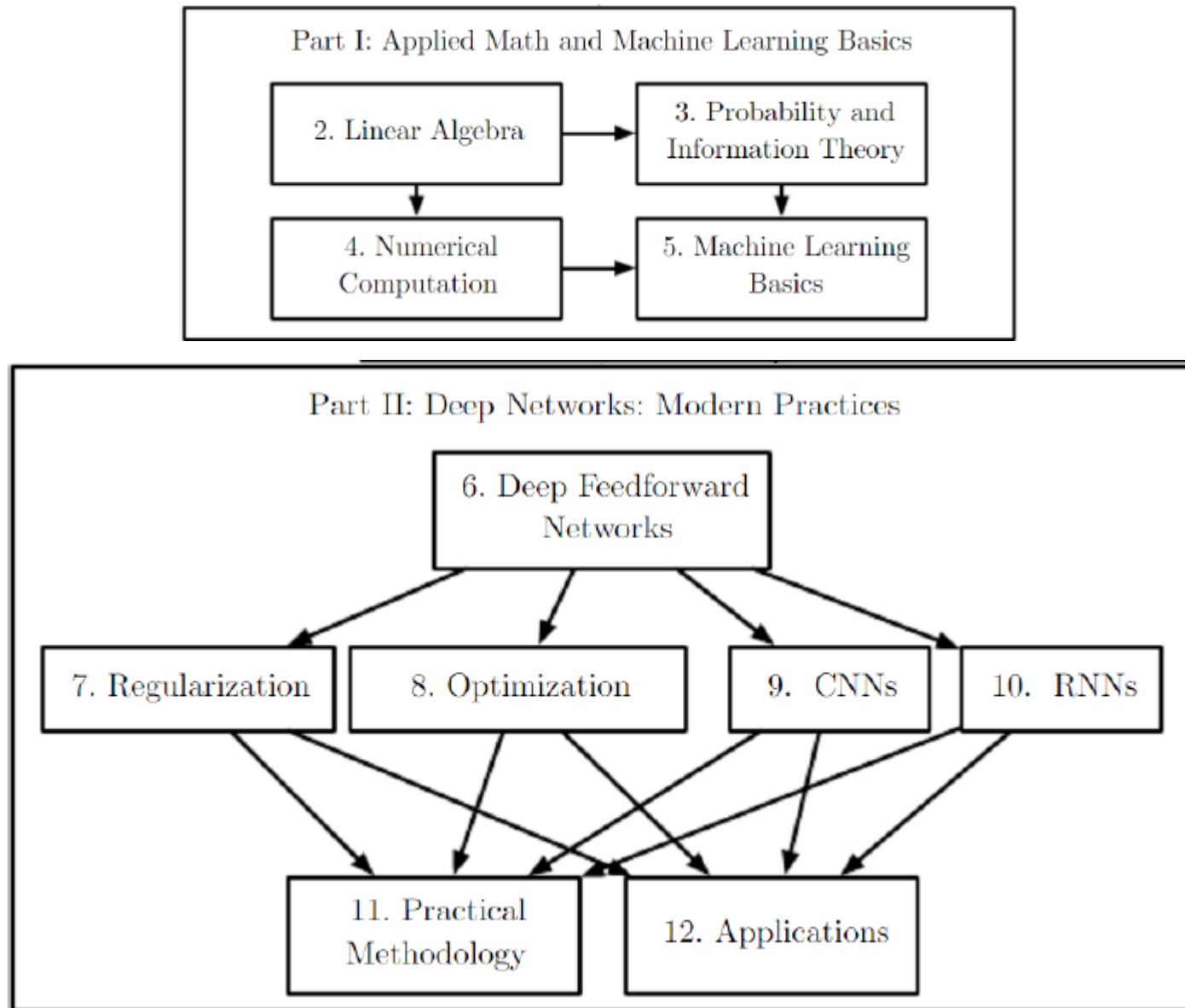


Deep Learning courses

Prof. Hung-yi Lee National Taiwan University (NTU) Taipei

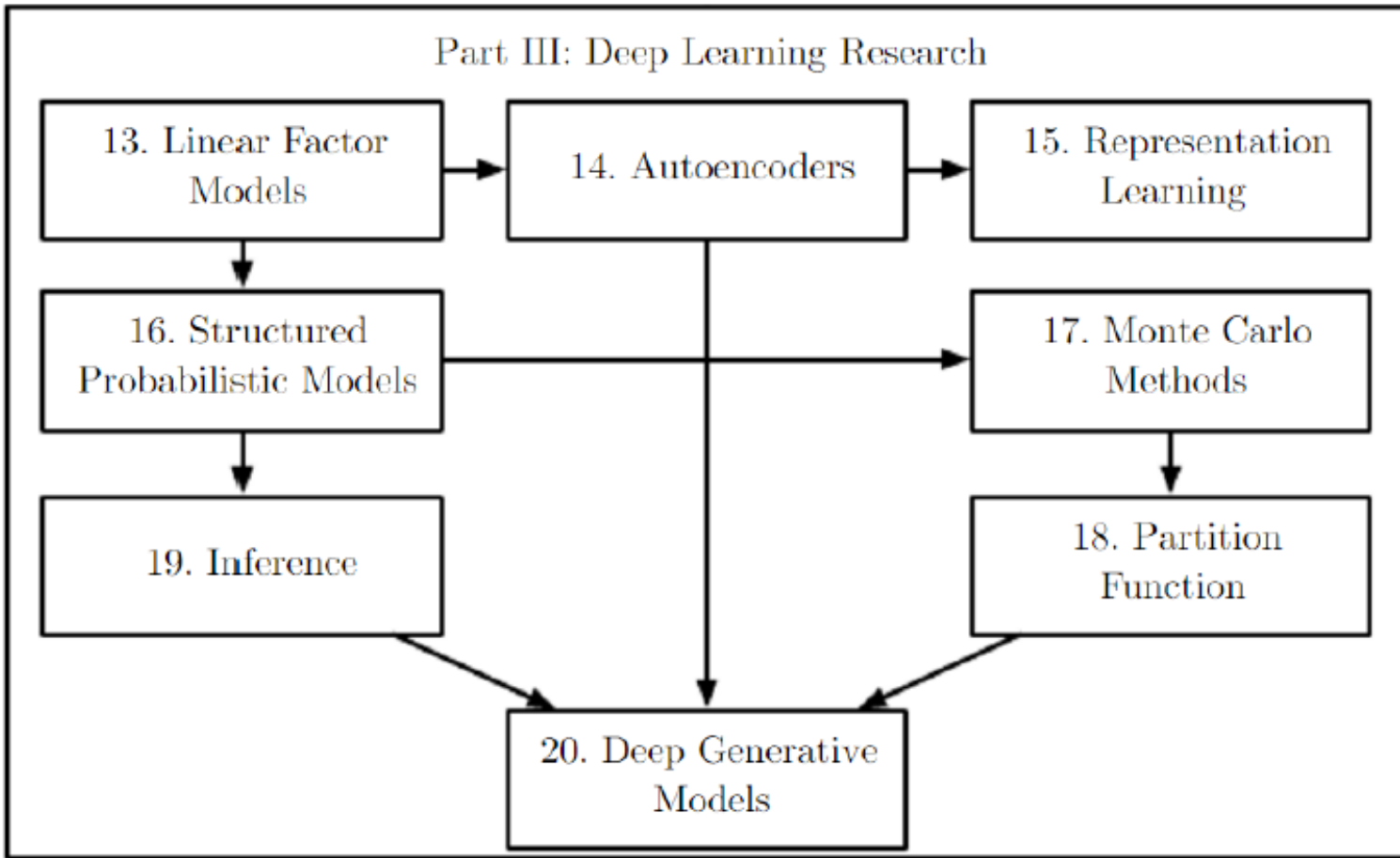
Introduction to Deep Learning

...the Machine Learning background...



Introduction to Deep Learning

...the Machine Learning background...



FROM : *humans*

- to *machines*
- to Artificial Intelligence (AI)

... so we will proceed

FROM:

- *manual classifiers*
- *linear classifiers*

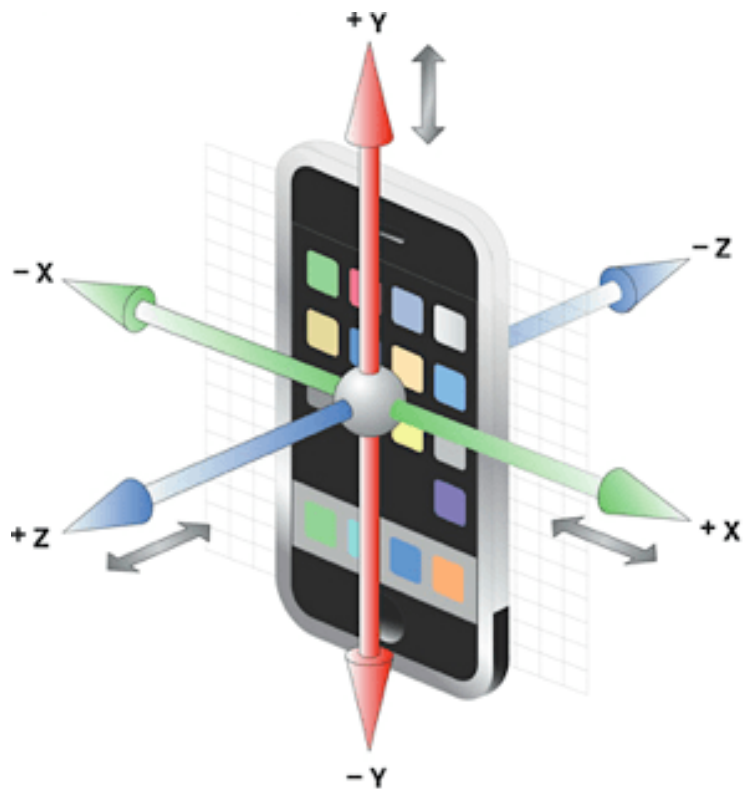
TO:

- *Neural Networks*

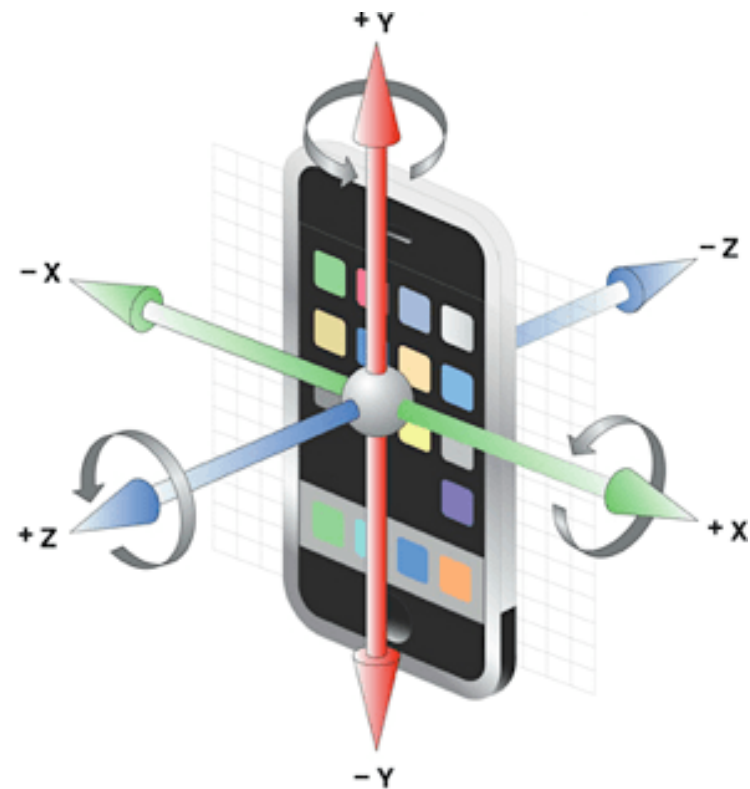
Drivies use case:

www.driviesapp.com

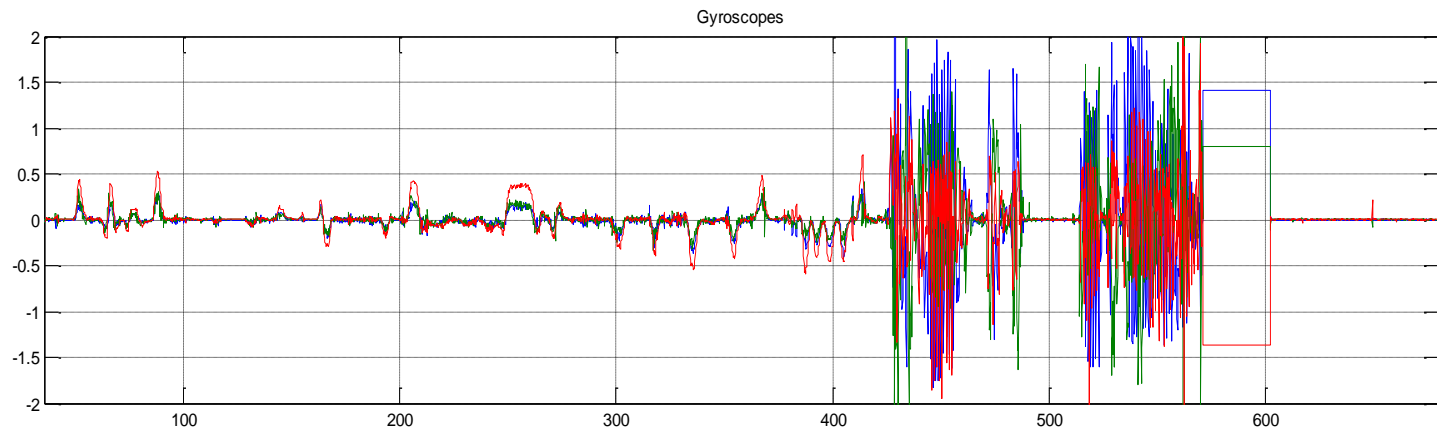
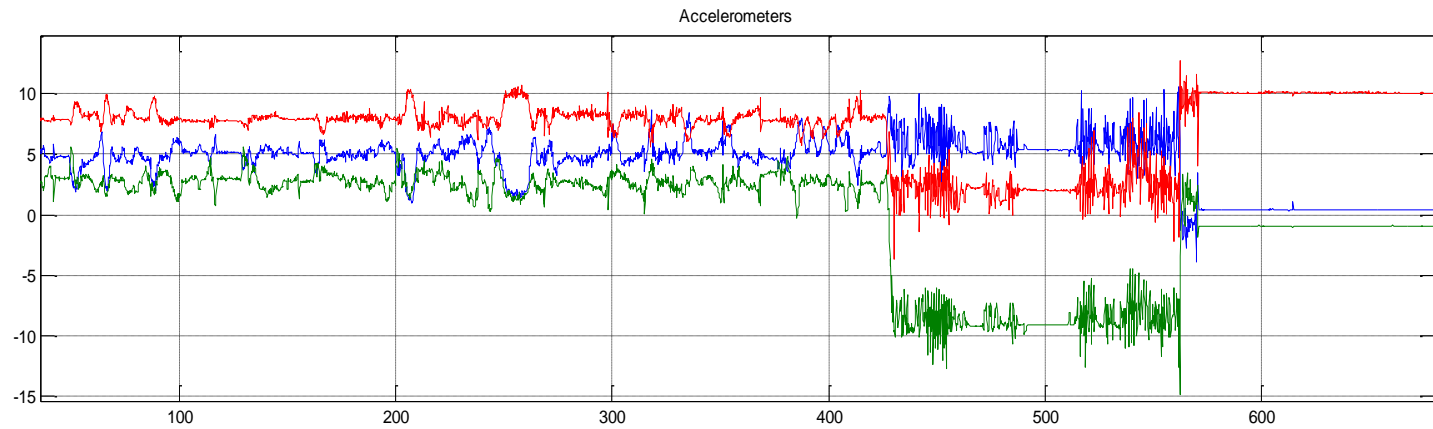




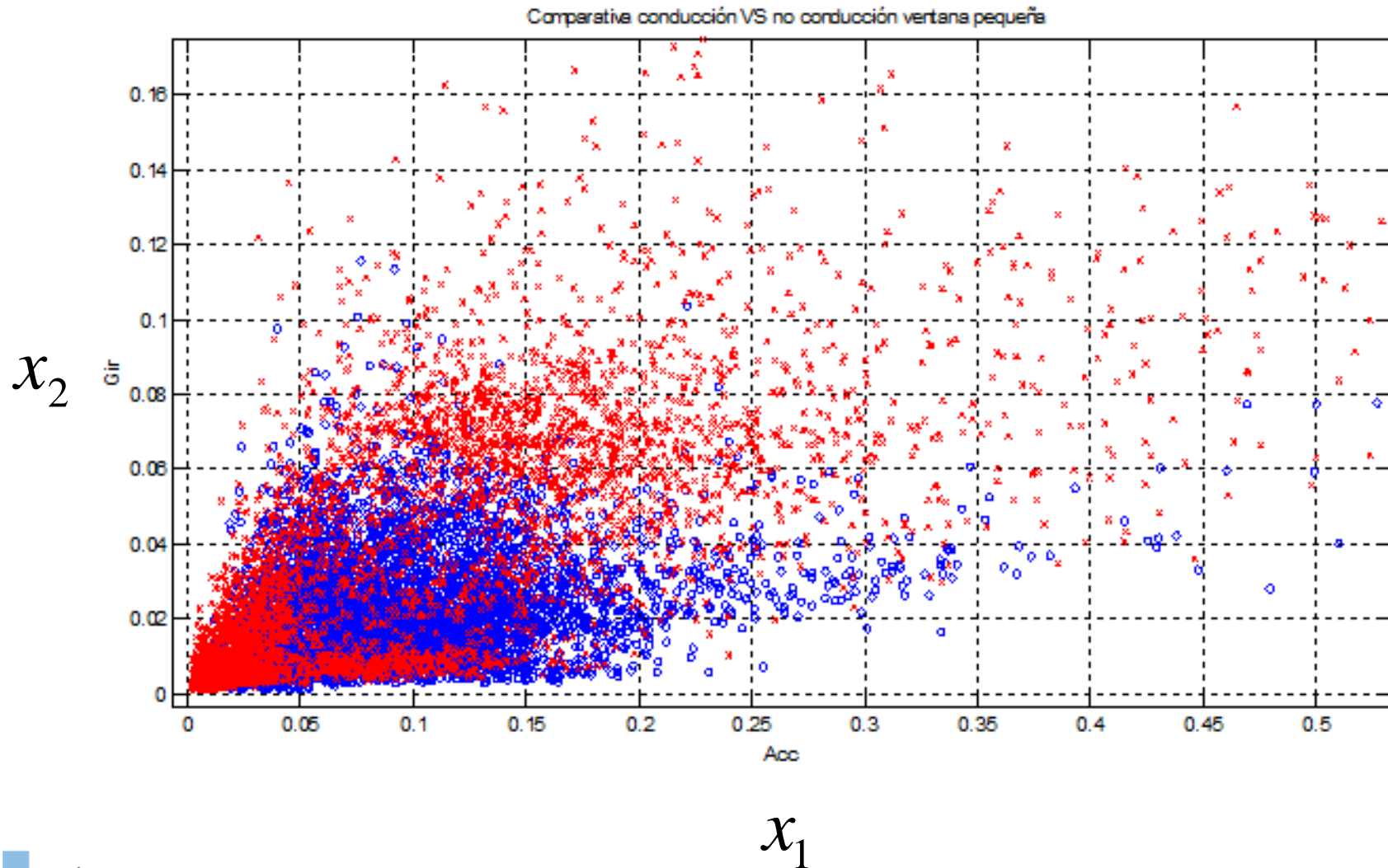
Accelerometer



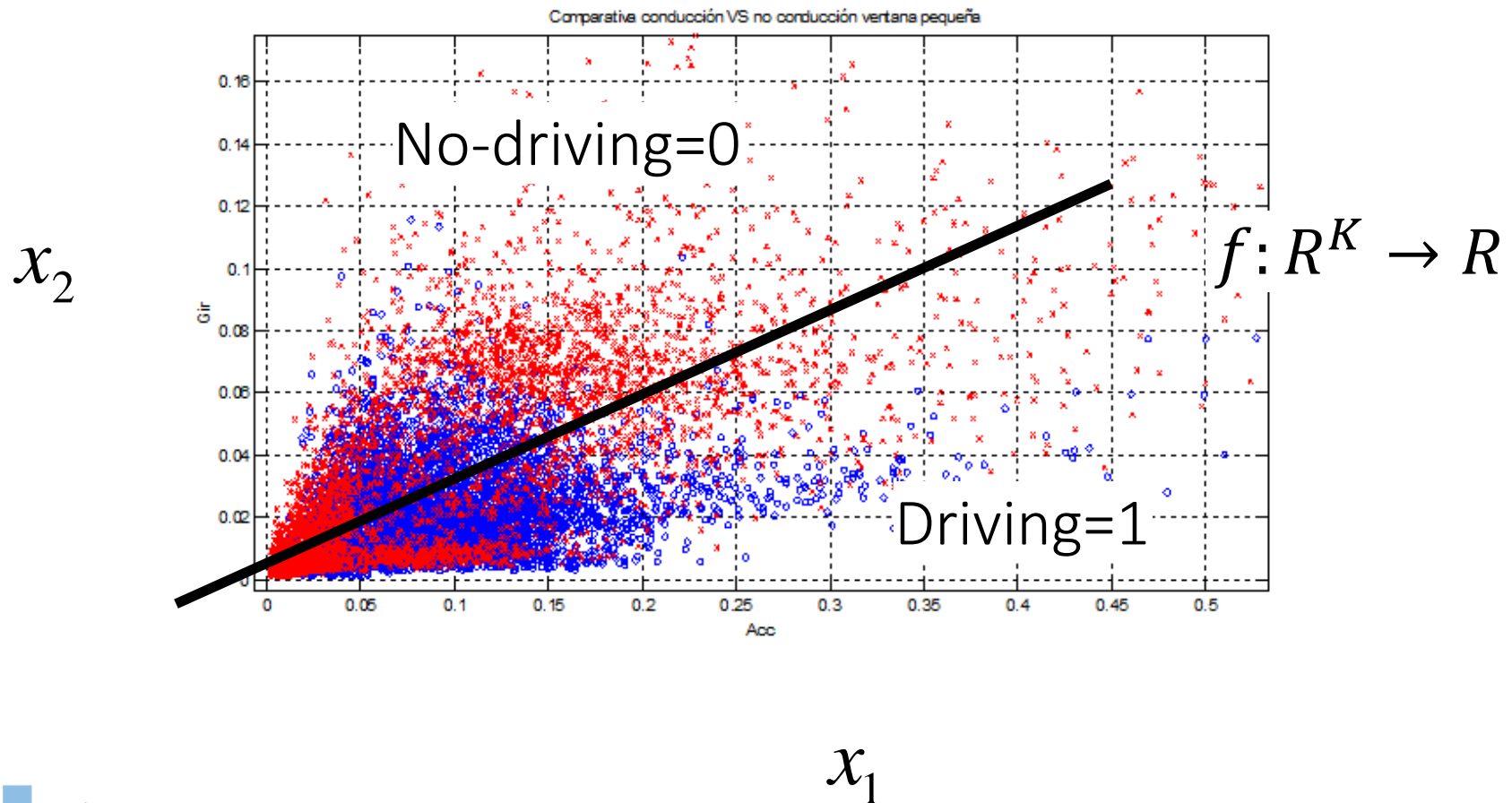
Gyroscope



From manual - linear classifiers TO Neural Networks

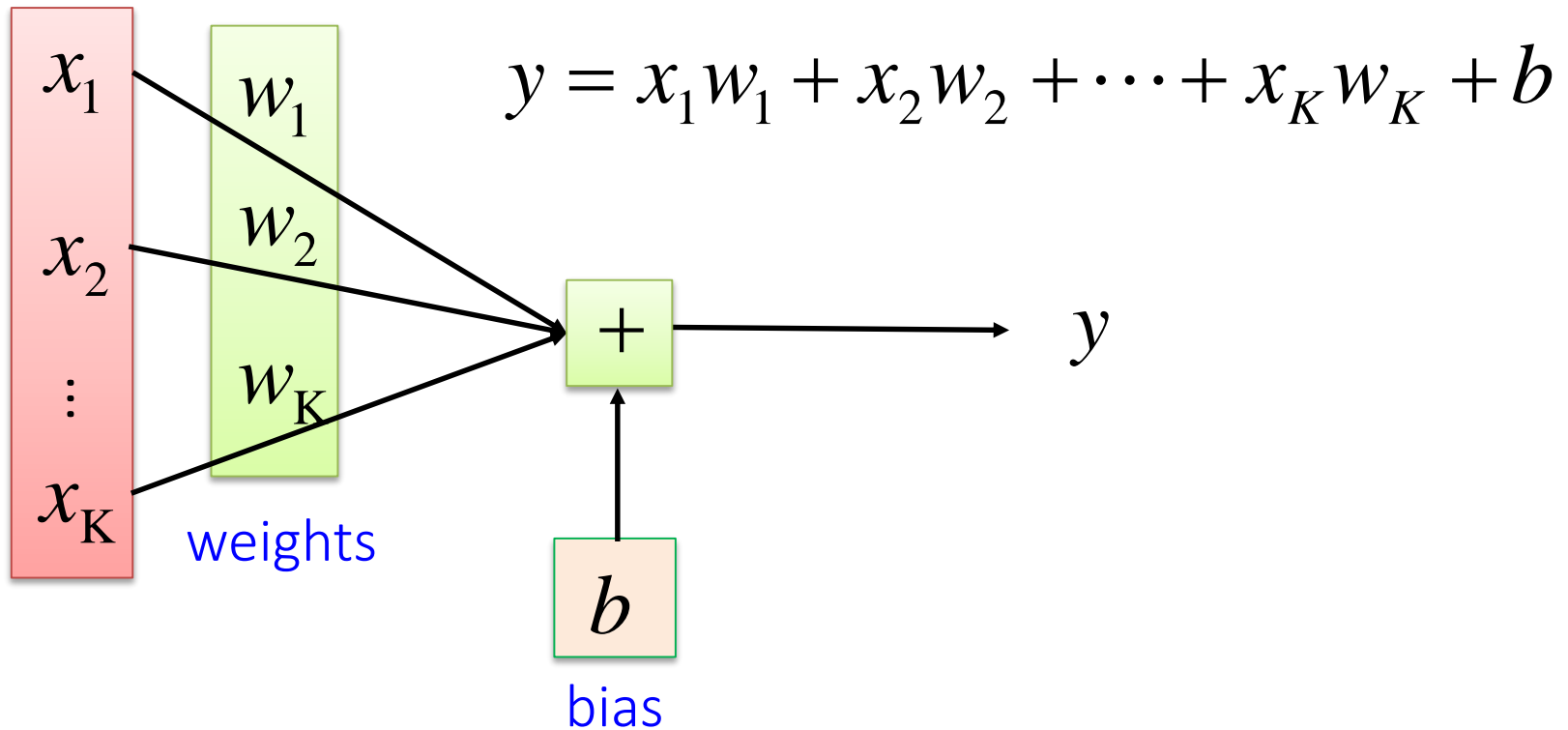


Driving detection (yes/no)
= define a **decision function**



From linear classifiers TO Neural Networks

A Linear decision function



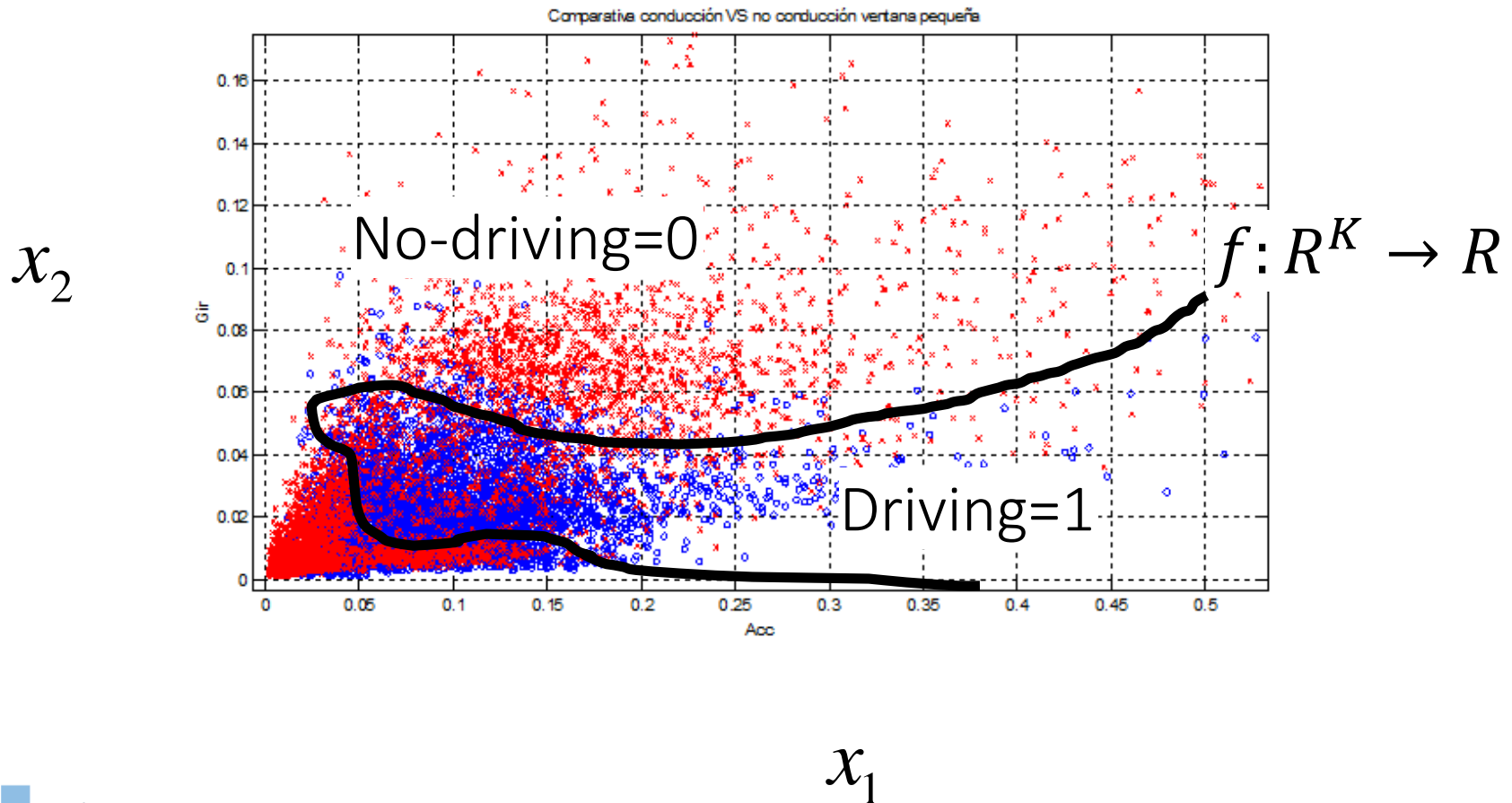
A Linear **decision** function

$$y = x_1 w_1 + x_2 w_2 + \cdots + x_K w_K + b$$

$$y = \mathbf{x}^T \mathbf{w}$$

From linear classifiers TO Neural Networks

Nonlinear decision function?



Non-linear decision functions

$$y = x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_1 x_2 w_4 + \dots + b$$

A linear model of transformed inputs:

$$y = \phi(\mathbf{x})^T \mathbf{w}$$

$\phi(\mathbf{x})$ where ϕ is a non linear transformation

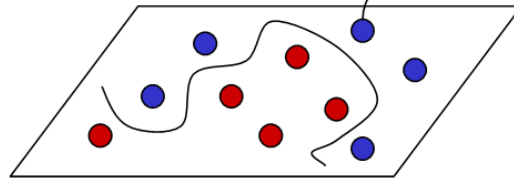
How choosing the mapping $\phi(.)$?

1. To manually engineer $\phi(.)$
2. Use a very generic $\phi(.)$ as kernel machines (e.g. SVM, RBF kernel)
3. The strategy of **deep learning** : to learn $\phi(.)$

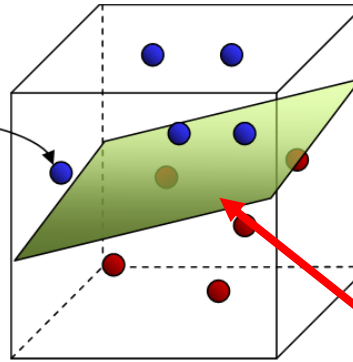
SVM

Hand-crafted
kernel function

ϕ



Input Space



Feature Space

Apply simple
classifier

From: DL Tutorial
Hung-yi Lee

Source of image: http://www.gipsa-lab.grenoble-inp.fr/transfert/seminaire/455_Kadri2013Gipsa-lab.pdf

The DL approach: learn $\phi(\cdot)$

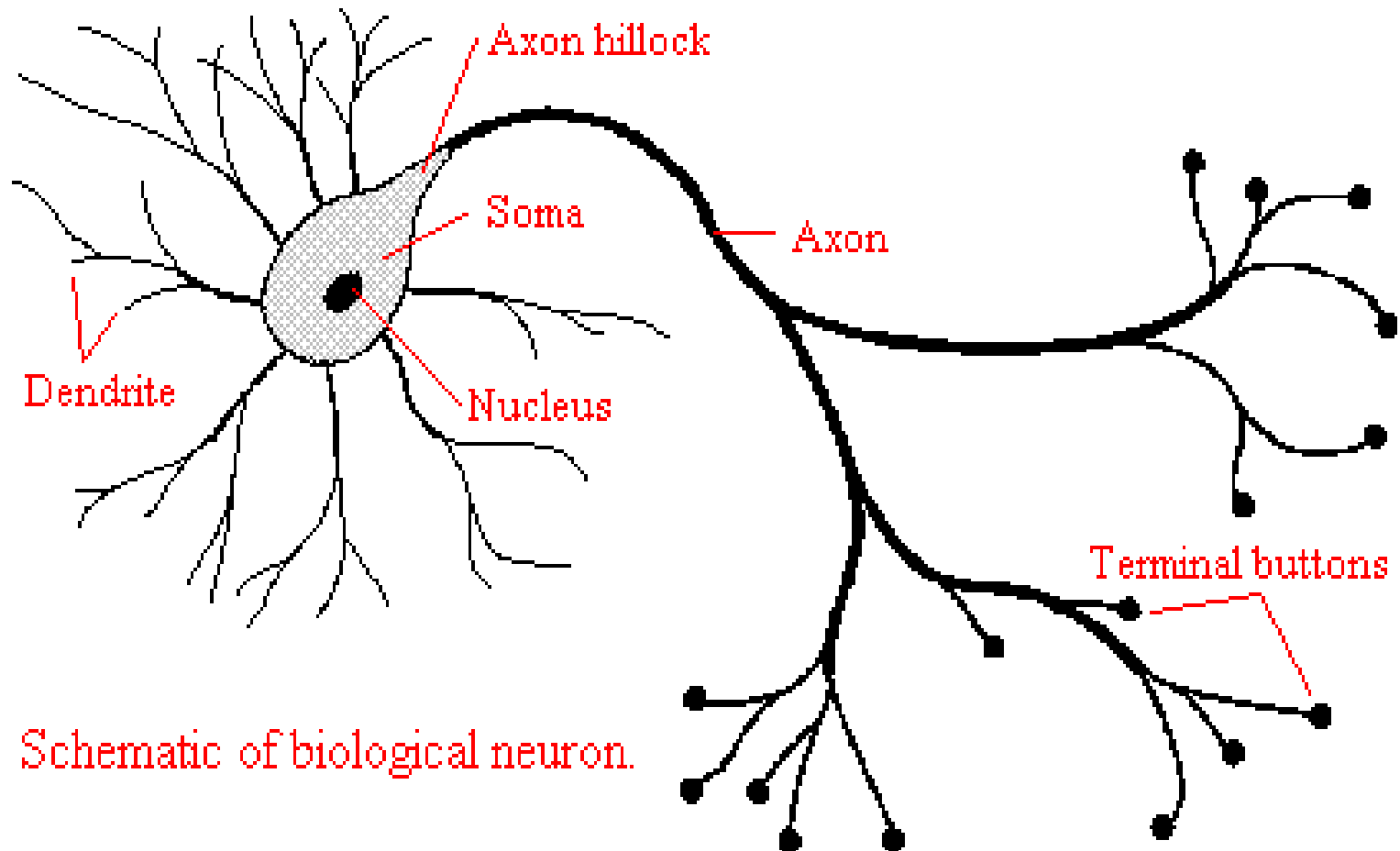
Now we have:

- Parameters $\boldsymbol{\theta}$ that we use to learn $\phi(\cdot)$ from a broad class of functions
- Parameters \mathbf{w} that map $\phi(\mathbf{x})$ to the desired output

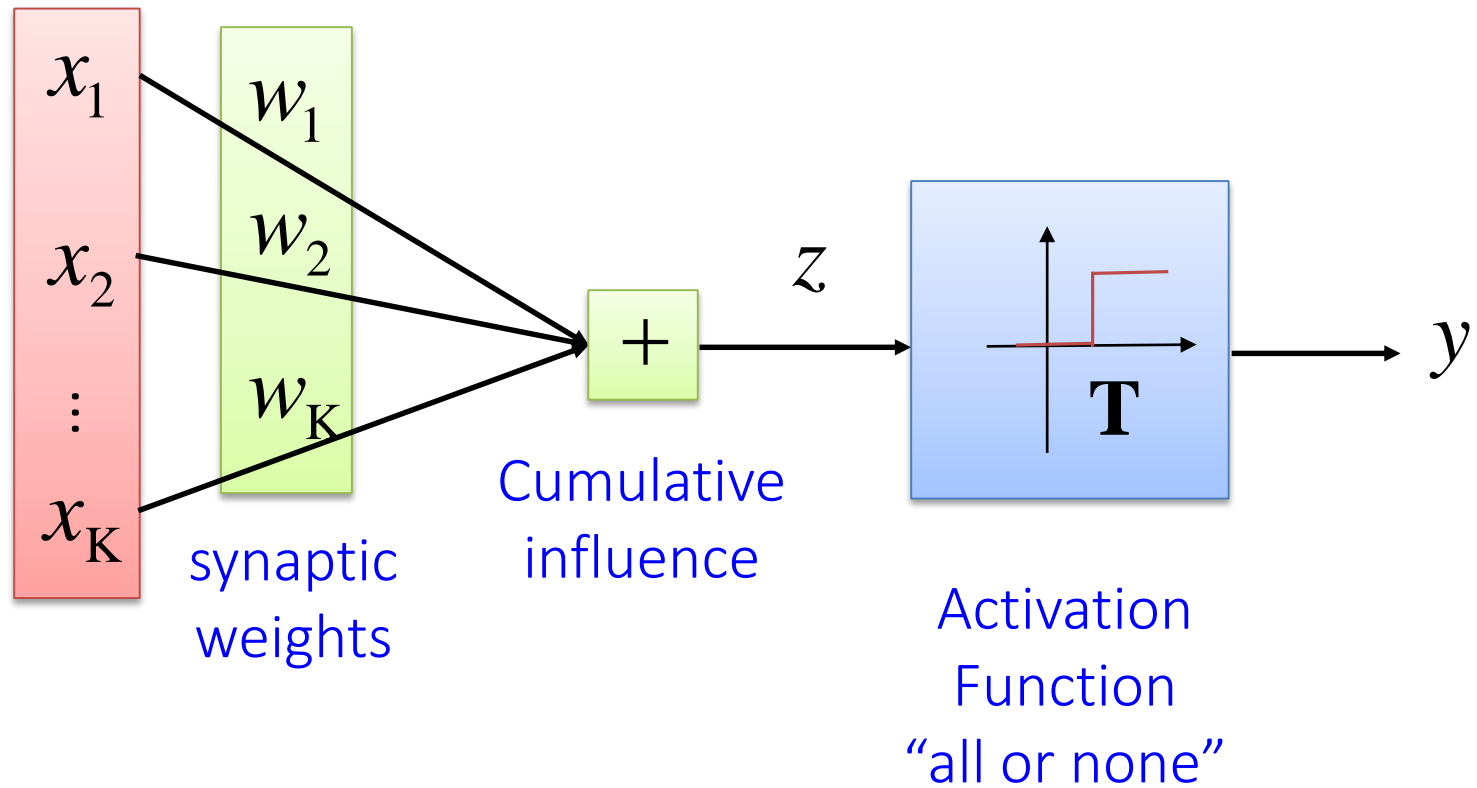
$$y = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \phi(\mathbf{x})^T \mathbf{w}$$

The DL approach: learn $\phi(\mathbf{x})$

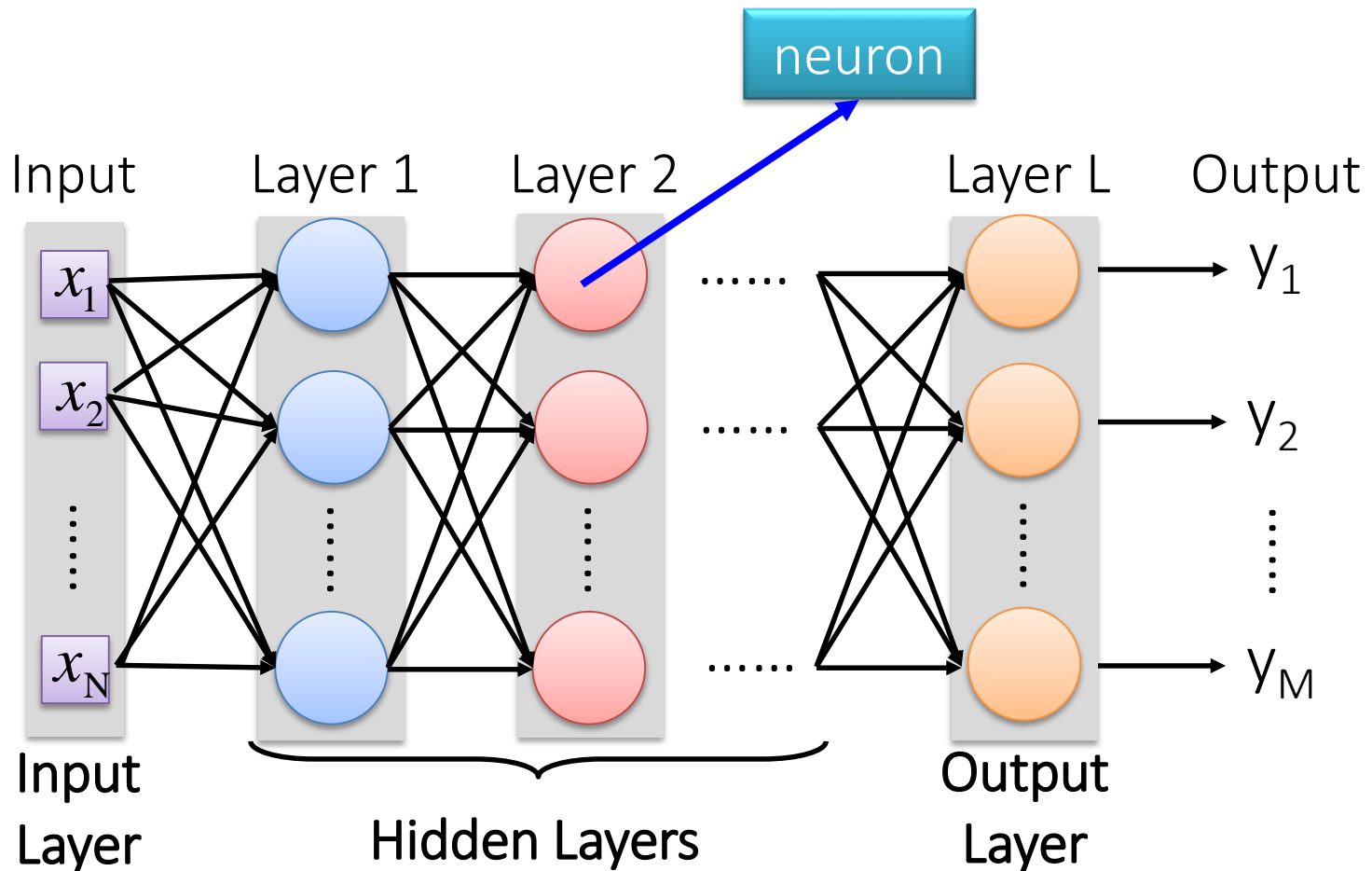
...from a broad class of functions



Neuron approach....



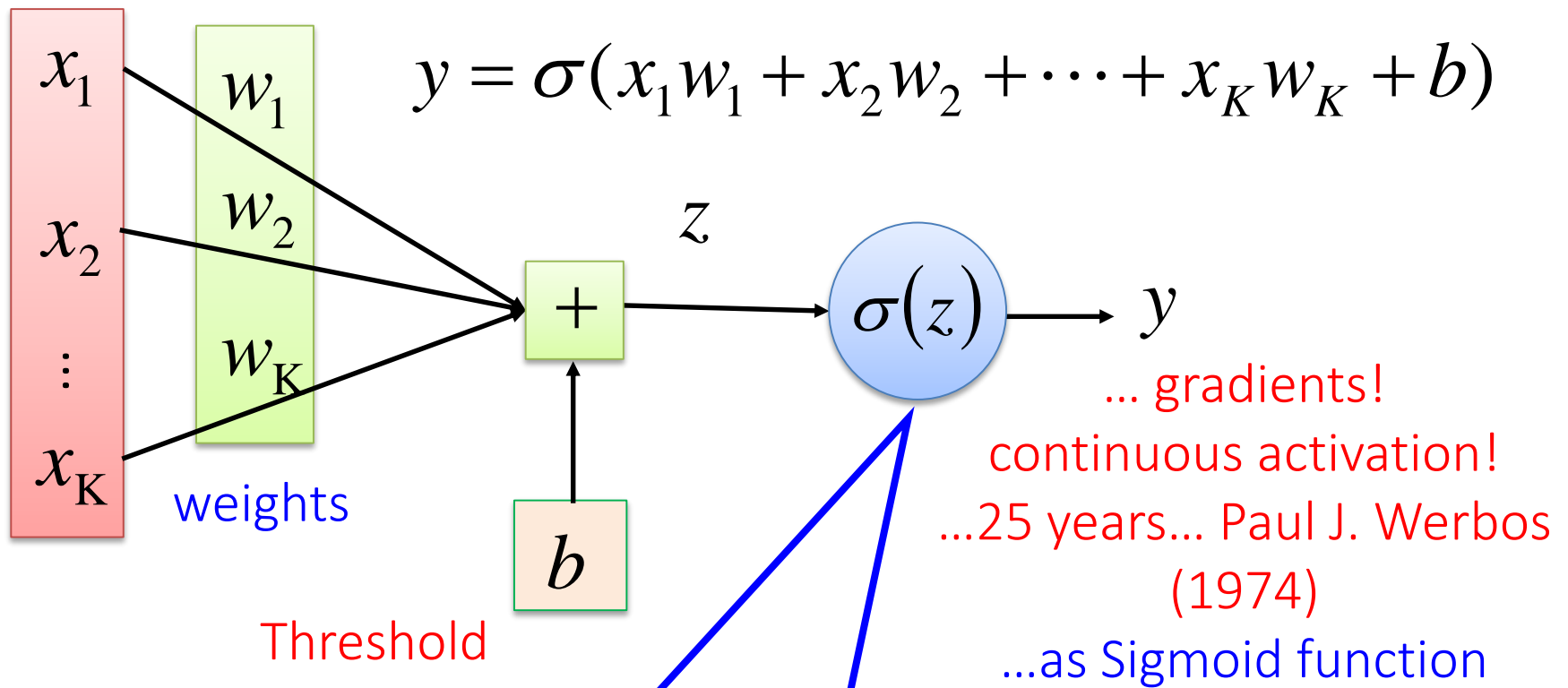
Neural Network (from Hung-yi Lee “Deep Learning Tutorial”)



Deep means many hidden layers

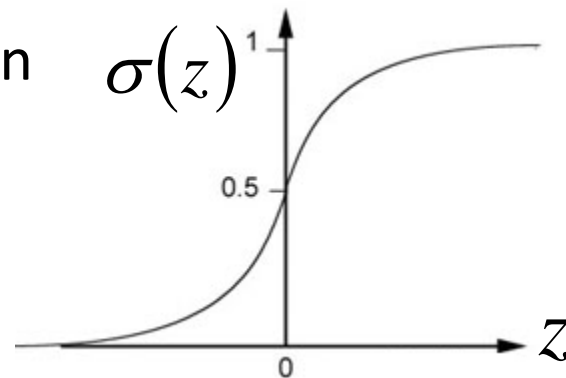
From linear classifiers TO Neural Networks

- Parameters θ that we use to learn $\phi(.)$ from a broad class of functions
- **Weights** and **thresholds** are estimated from training examples:
 - to minimize a **loss function** (i.e. similarity between NN outputs y and desired outputs \hat{y})

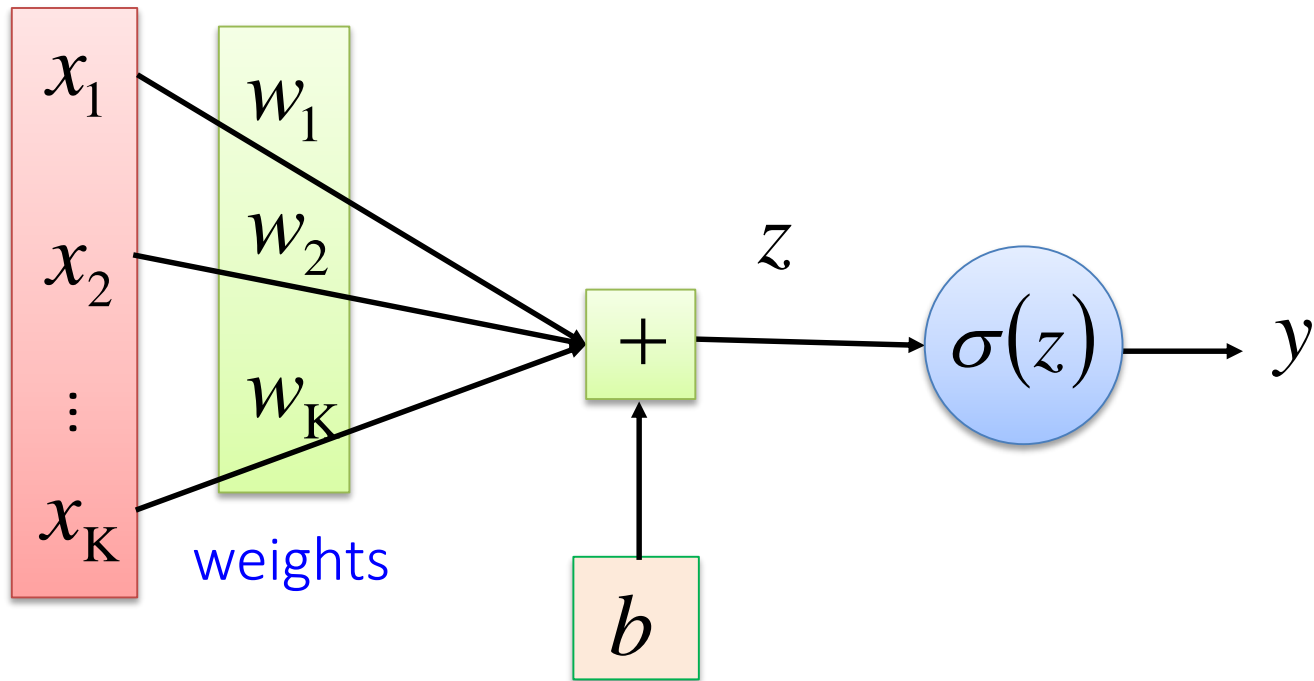


Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



- But recall that this is also logistic regression!



From linear classifiers TO Neural Networks

- So let's stop here and start playing with



Google TensorFlow



- Library for writing “machine intelligence” algorithms
- Very popular for deep learning and neural networks
- Can also be used for general purpose numerical computations
- Interface in C++ and Python

Deep learning libraries: growth over past three months

new contributors from 2016-10-09 to 2017-02-10

#1:	192	tensorflow/tensorflow
#2:	89	dmlc/mxnet
#3:	78	fchollet/keras
#4:	42	baidu/paddle
#5:	29	Microsoft/CNTK
#6:	23	pfnet/chainer
#7:	21	Theano/Theano
#8:	20	deeplearning4j/deeplearning4j
#9:	20	tflearn/tflearn
#10:	19	BVLC/caffe
#11:	9	torch/torch7
#12:	3	NVIDIA/DIGITS

new forks from 2016-10-09 to 2017-02-10

#1:	6525	tensorflow/tensorflow
#2:	1822	BVLC/caffe
#3:	1316	fchollet/keras
#4:	999	dmlc/mxnet
#5:	909	deeplearning4j/deeplearning4j
#6:	887	Microsoft/CNTK
#7:	324	tflearn/tflearn
#8:	321	baidu/paddle
#9:	287	Theano/Theano
#10:	257	torch/torch7
#11:	175	NVIDIA/DIGITS
#12:	142	pfnet/chainer

new issues from 2016-10-09 to 2017-02-10

#1:	1563	tensorflow/tensorflow
#2:	979	fchollet/keras
#3:	871	dmlc/mxnet
#4:	646	baidu/paddle
#5:	486	Microsoft/CNTK
#6:	361	deeplearning4j/deeplearning4j
#7:	318	BVLC/caffe
#8:	217	NVIDIA/DIGITS
#9:	214	Theano/Theano
#10:	167	tflearn/tflearn
#11:	150	pfnet/chainer
#12:	90	torch/torch7














aggregate metrics growth from 2016-10-09 to 2017-02-10

#1:	54.01	tensorflow/tensorflow
#2:	18.71	fchollet/keras
#3:	16.38	dmlc/mxnet
#4:	12.86	BVLC/caffe
#5:	10.17	Microsoft/CNTK
#6:	9.32	baidu/paddle
#7:	8.75	deeplearning4j/deeplearning4j
#8:	4.21	Theano/Theano
#9:	3.89	tflearn/tflearn
#10:	3.14	NVIDIA/DIGITS
#11:	2.90	pfnet/chainer
#12:	2.46	torch/torch7



Deep learning libraries: Accumulated GitHub metrics

Aggregate popularity $(30 \cdot \text{contrib} + 10 \cdot \text{issues} + 5 \cdot \text{forks}) \cdot 1e-3$

#1:	172.29		tensorflow/tensorflow
#2:	89.78		BVLC/caffe
#3:	69.70		fchollet/keras
#4:	53.09		dmlc/mxnet
#5:	38.23		Theano/Theano
#6:	29.86		deeplearning4j/deeplearning4j
#7:	27.99		Microsoft/CNTK
#8:	17.36		torch/torch7
#9:	14.43		baidu/paddle
#10:	13.10		pfnet/chainer
#11:	12.37		NVIDIA/DIGITS
#12:	10.42		tflearn/tflearn
#13:	9.20		pytorch/pytorch

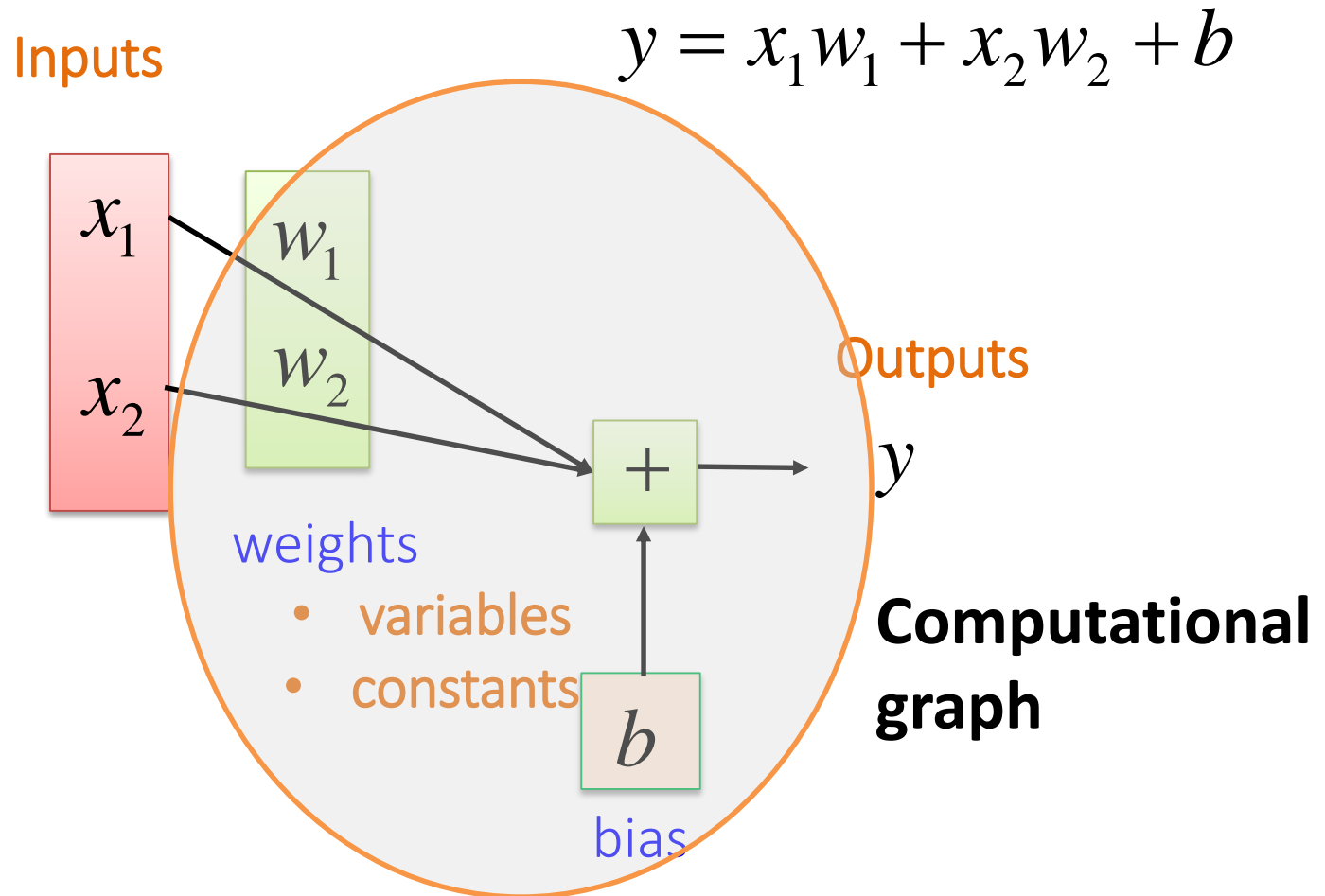




TensorFlow

Generates a computational graph like Theano
Everything about TensorFlow is here:

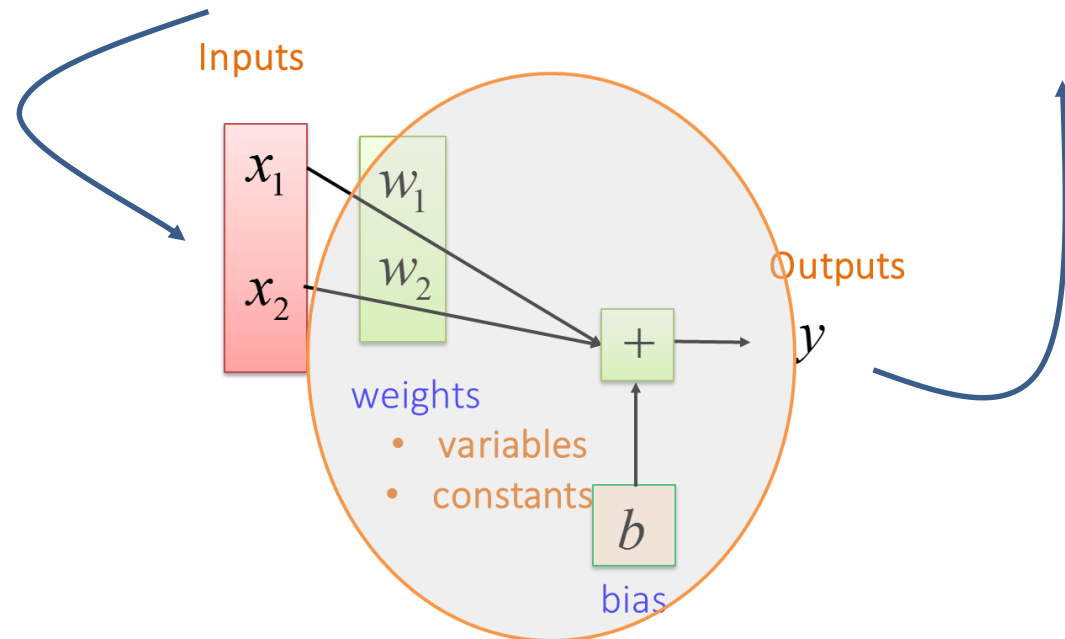
<https://www.tensorflow.org>



TensorFlow Recipe:



- Define a series of expressions
- Initialize variables
- Start a session (launch a graph)
- Run the graph, **feed** some data, **fetch** some values



TensorFlow Essentials:



Four types of objects make TensorFlow unique from other frameworks

- Session
- Computational graph
- Variables
- Placeholder

Let's start playing with



- Where? IBM DSX or your own (be careful with TF versions!!)
- How: we have prepared some Interactive Python notebooks (Jupyter) <http://jupyter.org/about.html>

MSTC_IntroTF_1.ipynb



MSTC_IntroTF_1.ipynb

The screenshot shows a Jupyter Notebook interface. At the top, there's a dark blue header with the 'Data Science Experience' logo and navigation links: 'Projects', 'Tools', 'Data Services', and 'Community'. Below this is a breadcrumb trail: 'My Projects > MSTC_DeepLearning > MSTC_IntroTF_1'. The main toolbar includes 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', and 'Help'. A secondary toolbar contains icons for saving, adding, deleting, and running cells, along with a 'Format' dropdown set to 'Markdown' and a 'CellToolbar' button. The notebook content area displays a bold heading 'This notebook contains a simple introduction to TF' followed by 'as used in:'. A bulleted list contains 'MSTC seminar on Deep Learning & Tensorflow'. Below this is a code cell with the following Python code:

```
In [1]: import tensorflow as tf

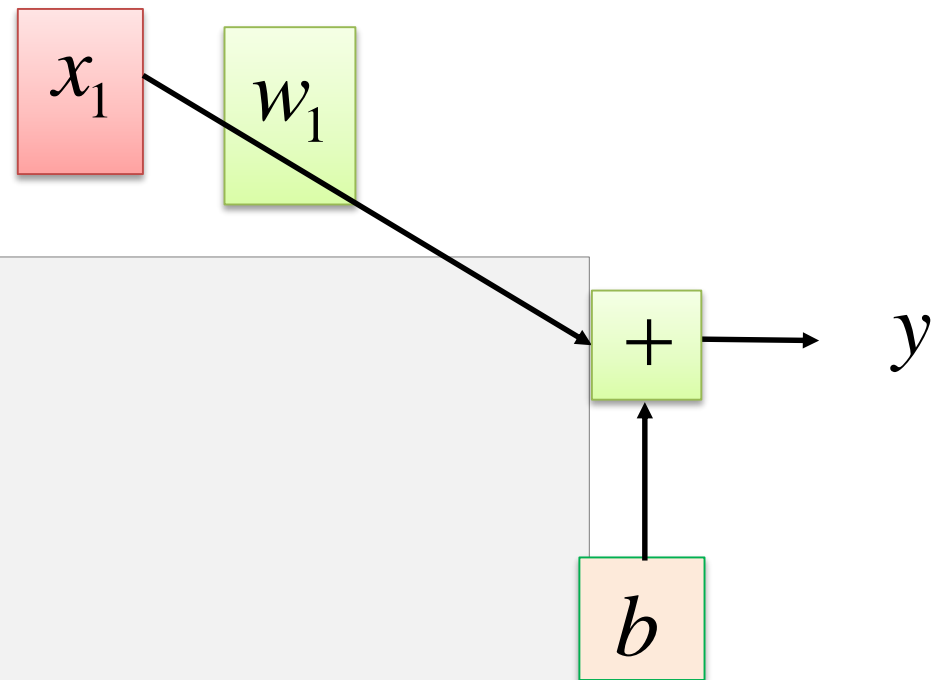
x=tf.constant(1.0)
W=tf.constant(6.0)
b=tf.constant(1.5)

y=x*W+b
```

TensorFlow Docs:

“TensorFlow programs are usually structured into a construction phase, that assembles a graph...”

$$y = x_1 w_1 + b$$



```
import tensorflow as tf

x=tf.constant(1.0)
w=tf.constant(6.0)
b=tf.constant(1.5)

y=x*w+b

print(y)
Tensor("add:0", shape=(), dtype=float32)
```

TensorFlow Docs:

- “A Session object encapsulates the environment in which Tensor objects are evaluated...”
- “..and an execution phase that uses a session to execute ops in the graph”

```
import tensorflow as tf

x=tf.constant(1.0)
w=tf.constant(6.0)
b=tf.constant(1.5)

y=x*w+b

with tf.Session() as sess:
    print(sess.run(y))
```

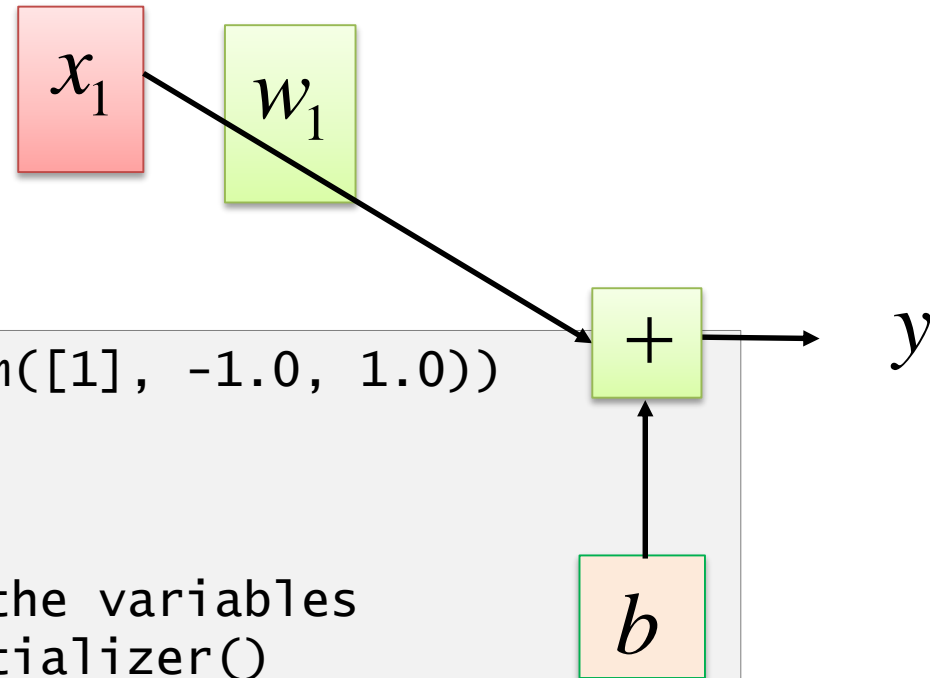
7.5

TensorFlow Variables

- “... hold and update **parameters**”

$$y = x_1 w_1 + b$$

- ..and graph (session)
is executed several times



```
w=tf.Variable(tf.random_uniform([1], -1.0, 1.0))
b=tf.Variable(tf.zeros([1]))
x=tf.constant(1.0)
```

```
# Before starting, initialize the variables
init = tf.global_variables_initializer()
```

```
y=x*w+b
```

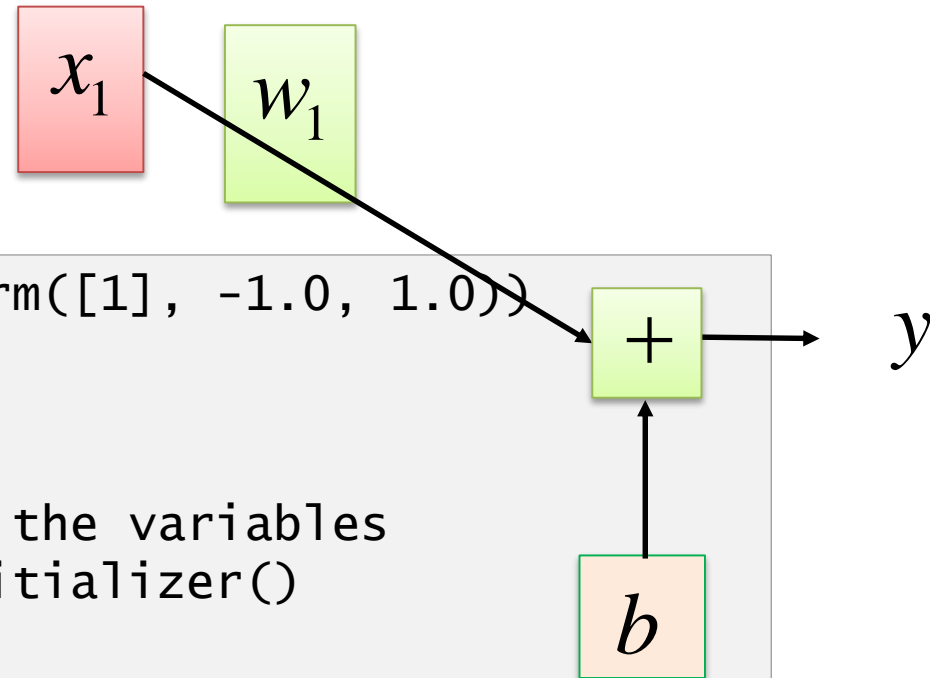
```
with tf.Session() as sess:
    sess.run(init)
    for step in range(4):
        print(sess.run(y))
```

TensorFlow Placeholders



- “... dummy nodes that provide **entry points** to the computational graph”

$$y = x_1 w_1 + b$$



```
W=tf.Variable(tf.random_uniform([1], -1.0, 1.0))
b=tf.Variable(tf.zeros([1]))
x=tf.placeholder(tf.float32)
```

```
# Before starting, initialize the variables
init = tf.global_variables_initializer()
```

```
y=x*w+b
```

```
with tf.Session() as sess:
    sess.run(init)
    for step in range(4):
        print(sess.run(y, feed_dict=(x:1)))
```

Why TensorFlow?



- Python + Numpy
- Graph based, easy to model
- Faster compile times than Theano
- **Tensorboard** for Visualization
- Open Sourced
- Data and Model Parallelism
- Distributed supported

But what's a Tensor?



Tensor Ranks, Shapes, and Types

Briefly: *A tensor is an array of n -dimension containing the same type of data*

- **Tensor rank** (sometimes referred to as order or degree or n -dimension) is the number of dimensions of the tensor.

Rank	Math entity	Python example
0	Scalar (magnitude only)	<code>s = 483</code>
1	Vector (magnitude and direction)	<code>v = [1.1, 2.2, 3.3]</code>
2	Matrix (table of numbers)	<code>m = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]</code>
3	3-Tensor (cube of numbers)	<code>t = [[[2], [4], [6]], [[8], [10], [12]], [[14], [16], [18]]]</code>
n	n -Tensor (you get the idea)	<code>....</code>

But what's a Tensor?



Tensor Ranks, Shapes, and Types

- **Tensor shape:** The TensorFlow documentation uses three notational conventions to describe tensor dimensionality: rank, shape, and dimension number.

Rank	Shape	Dimension number	Example
0	[]	0-D	A 0-D tensor. A scalar.
1	[D0]	1-D	A 1-D tensor with shape [5].
2	[D0, D1]	2-D	A 2-D tensor with shape [3, 4].
3	[D0, D1, D2]	3-D	A 3-D tensor with shape [1, 4, 3].
n	[D0, D1, ... Dn-1]	n-D	A tensor with shape [D0, D1, ... Dn-1].

Shapes can be represented via Python lists / tuples of ints, or with the `tf.TensorShape`.

But what's a Tensor?



Tensor Ranks, Shapes, and Types

Data type	Python type	Description
DT_FLOAT	<code>tf.float32</code>	32 bits floating point.
DT_DOUBLE	<code>tf.float64</code>	64 bits floating point.
DT_INT8	<code>tf.int8</code>	8 bits signed integer.
DT_INT16	<code>tf.int16</code>	16 bits signed integer.
DT_INT32	<code>tf.int32</code>	32 bits signed integer.
DT_INT64	<code>tf.int64</code>	64 bits signed integer.
DT_UINT8	<code>tf.uint8</code>	8 bits unsigned integer.
DT_UINT16	<code>tf.uint16</code>	16 bits unsigned integer.
DT_STRING	<code>tf.string</code>	Variable length byte arrays. Each element of a Tensor is a byte array.
DT_BOOL	<code>tf.bool</code>	Boolean.

.....

.....

But what's a Tensor?



A Tensor in TensorFlow has 2 shapes! The **static shape**
AND the **dynamic shape**

- The static shape can be read using the `tf.Tensor.get_shape()` method: this shape is inferred from the operations that were used to create the tensor, and may be partially complete.

```
x = tf.placeholder(tf.int32, shape=[4])  
print x.get_shape()  
# ==> '(4,)'
```

- If the static shape is not fully defined, the dynamic shape of a Tensor `t` can be determined by evaluating `tf.shape(t)`.

```
y, _ = tf.unique(x)  
print y.get_shape()  
# ==> '(?,)'
```

But what's a Tensor?



A Tensor in TensorFlow has 2 shapes! The **static shape**
AND the **dynamic shape**

- getting dynamic shape of a Tensor `t` can be determined by evaluating `tf.shape(t)`.

```
y, _ = tf.unique(x)
```

```
sess = tf.Session()  
print sess.run(y, feed_dict={x: [0, 1, 2, 3]}).shape  
# ==> '(4,)'
```

```
print sess.run(y, feed_dict={x: [0, 0, 0, 0]}).shape  
# ==> '(1,)'
```


But what's a Tensor?



mathematical operations to manipulate the *tensors*

Operation	Description
<code>tf.add</code>	sum
<code>tf.sub</code>	subtraction
<code>tf.mul</code>	multiplication
<code>tf.div</code>	division
<code>tf.mod</code>	module
<code>tf.abs</code>	return the absolute value
<code>tf.neg</code>	return negative value
<code>tf.sign</code>	return the sign
<code>tf.inv</code>	returns the inverse
<code>tf.square</code>	calculates the square
<code>tf.round</code>	returns the nearest integer
<code>tf.sqrt</code>	calculates the square root
<code>tf.pow</code>	calculates the power
<code>tf.exp</code>	calculates the exponential
<code>tf.log</code>	calculates the logarithm

.....

But what's a Tensor?



Dealing with shapes is one of the major issues when working with TensorFlow!

Placeholders: Feeding data...



`tf.placeholder()`

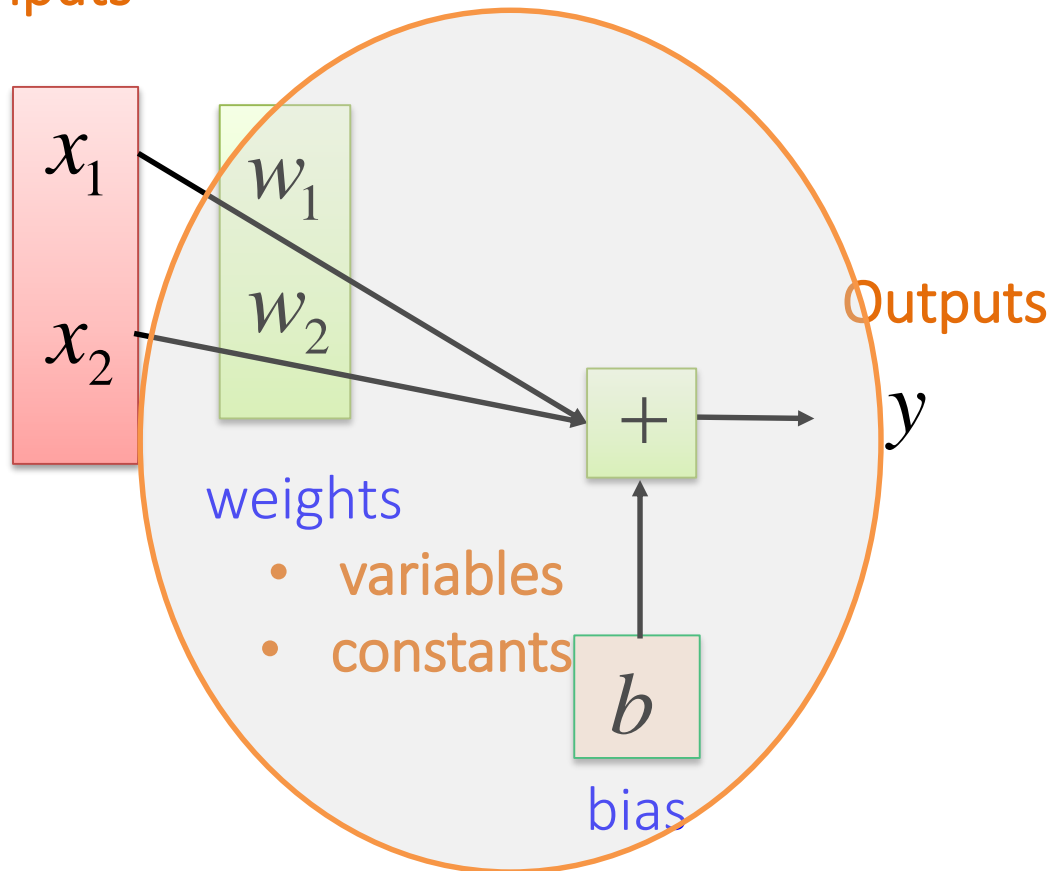
- Arguments: data type (mandatory), shape (optional), name (optional)
- Specifying (the entirety or part of) shape restricts shapes of possible inputs to the node
- At run time, we feed using a `feed_dict`

Placeholders allow you to pass in numpy arrays of data

Let's try a simple linear classifier

$$y = x_1 w_1 + x_2 w_2 + b$$

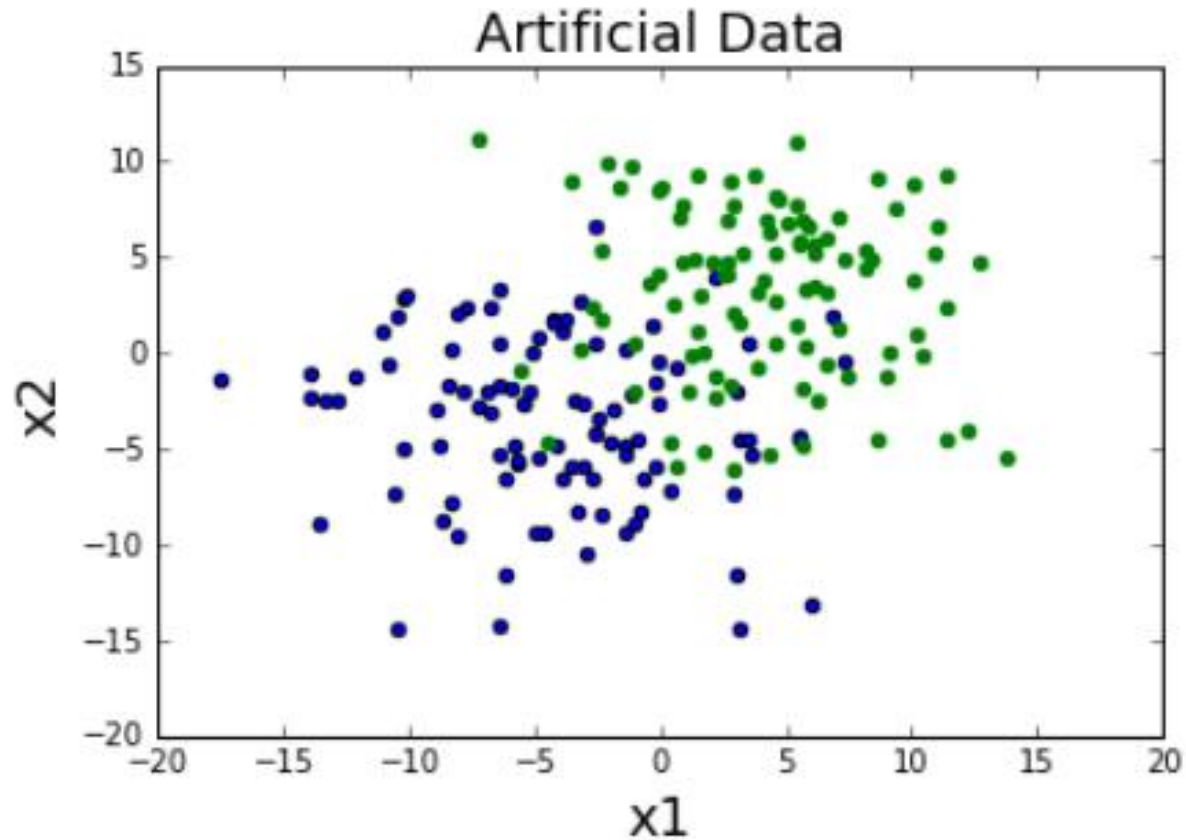
Inputs





TensorFlow

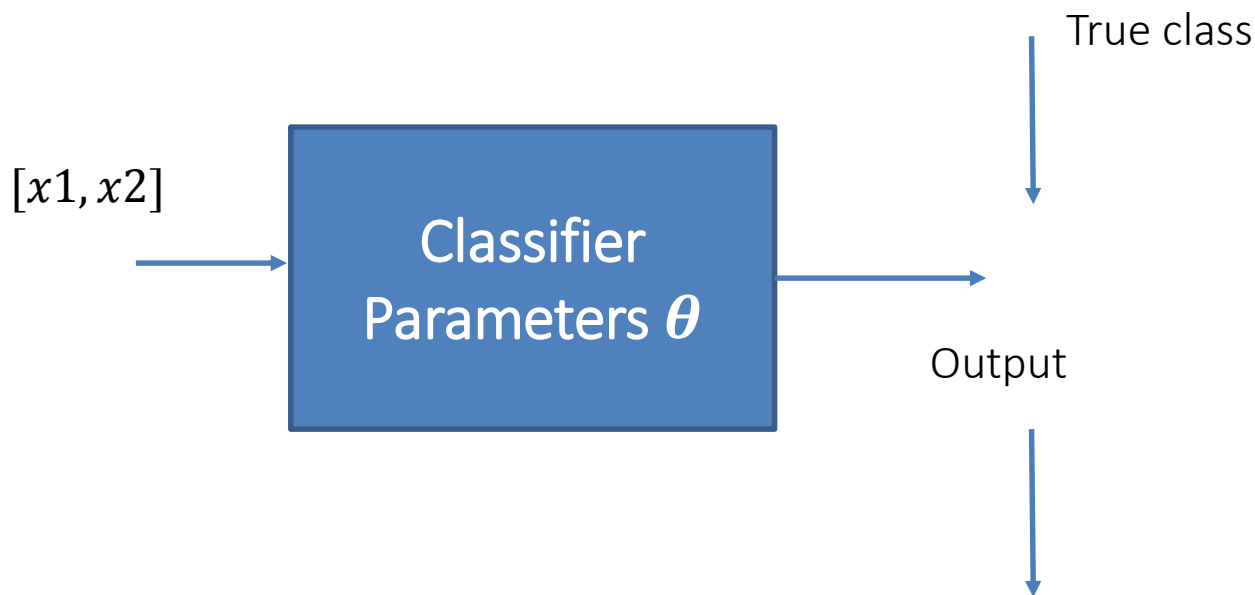
Binary (two-class) classifier from synthetic data





Binary (two-class) classifier from synthetic data

Training: How to find the parameters $\theta : W, b$?



OPTIMIZATION =>

Classification Error
Cost a function of the parameters = $C(\theta)$

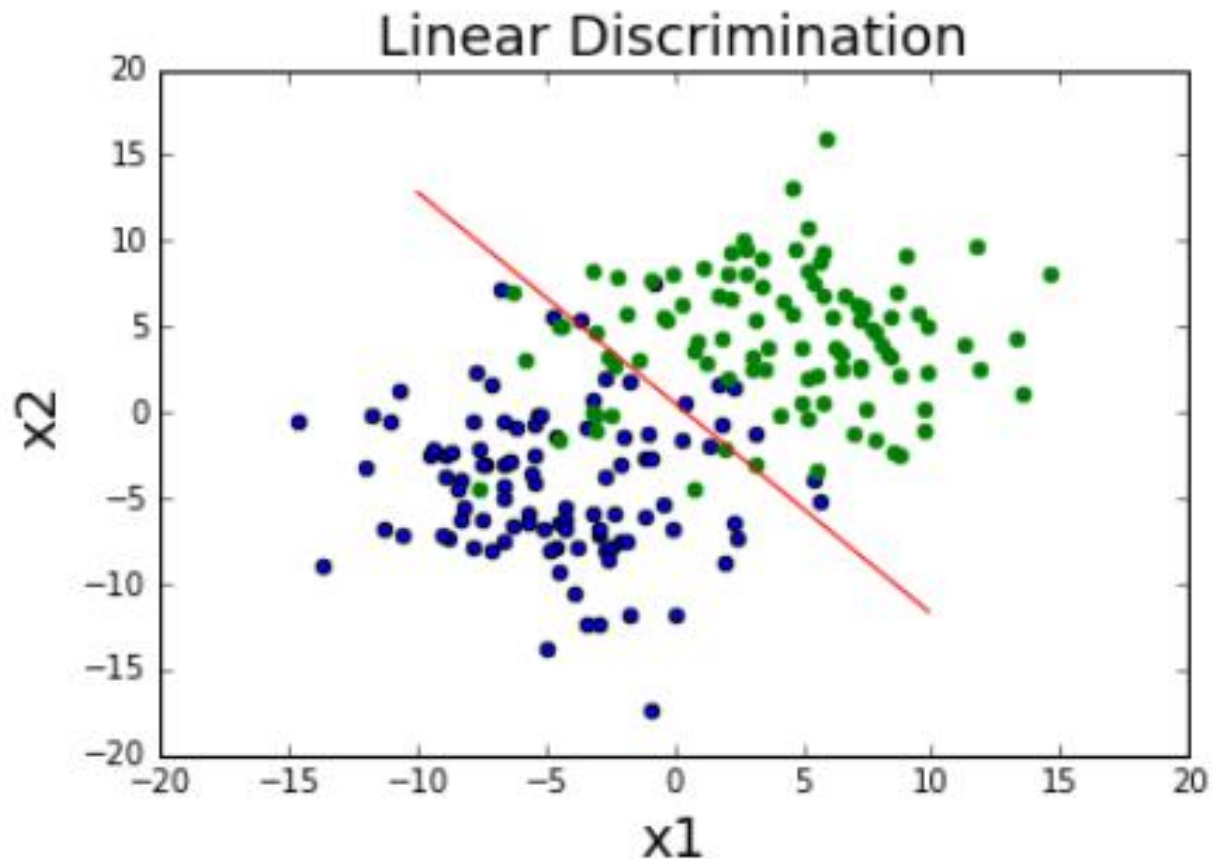
MSTC_IntroTF_2.ipynb



For particular values of W and b :

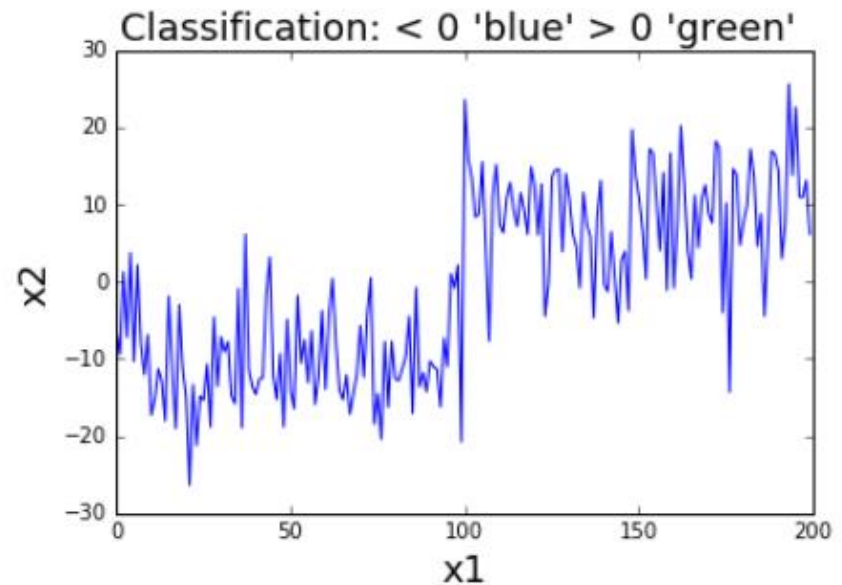
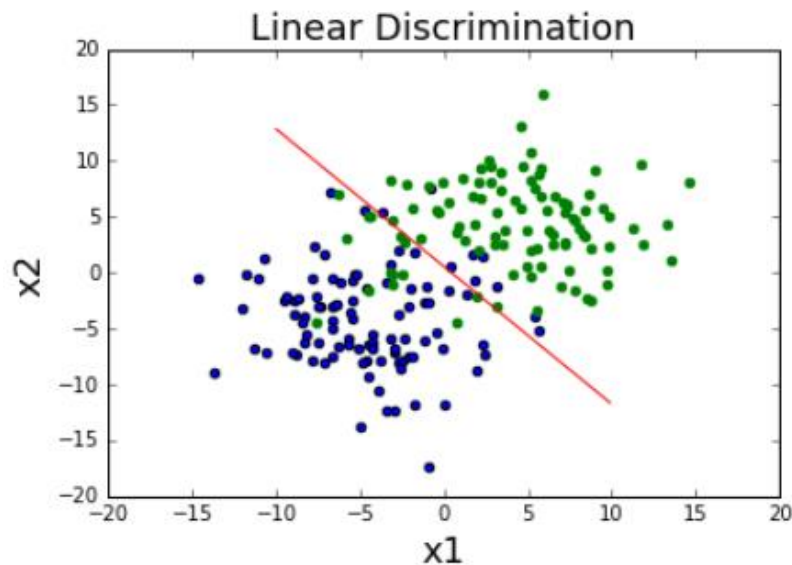
Red line is the set of points defined by: $x_1 w_1 + x_2 w_2 + b = 0$

... arbitrary $x_2 + 1.23 * x_1 - 0.55 = 0$



Green points will give: $x_2 + 1.23 * x_1 - 0.55 > 0$

Blue points will give: $x_2 + 1.23 * x_1 - 0.55 < 0$



Next Slides are from:
Deep Learning Tutorial

李宏毅

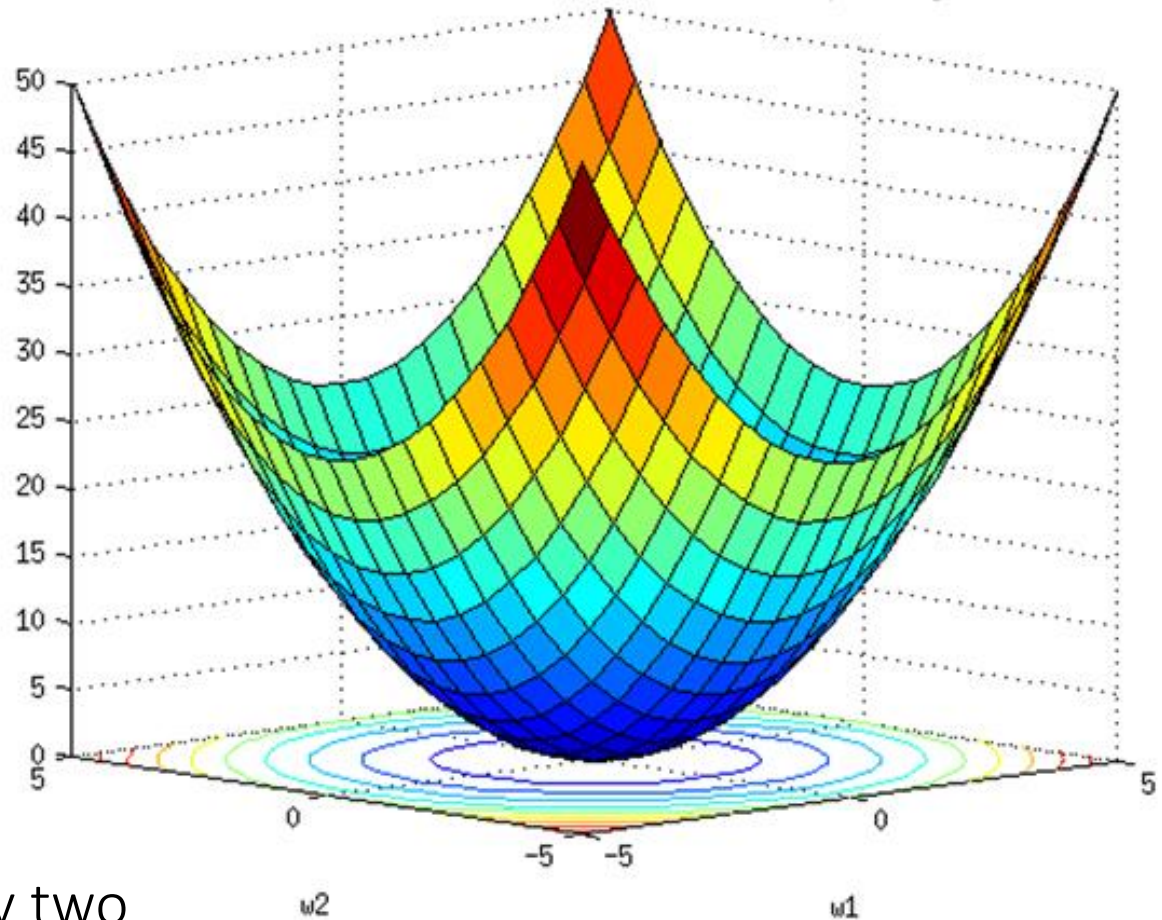
Hung-yi Lee

OPTIMIZATION

Gradient Descent

Cost
Or
Loss function

$$C(\theta)$$



Assume there are only two parameters w_1 and w_2 in a network.

$$\theta = \{w_1, w_2\}$$

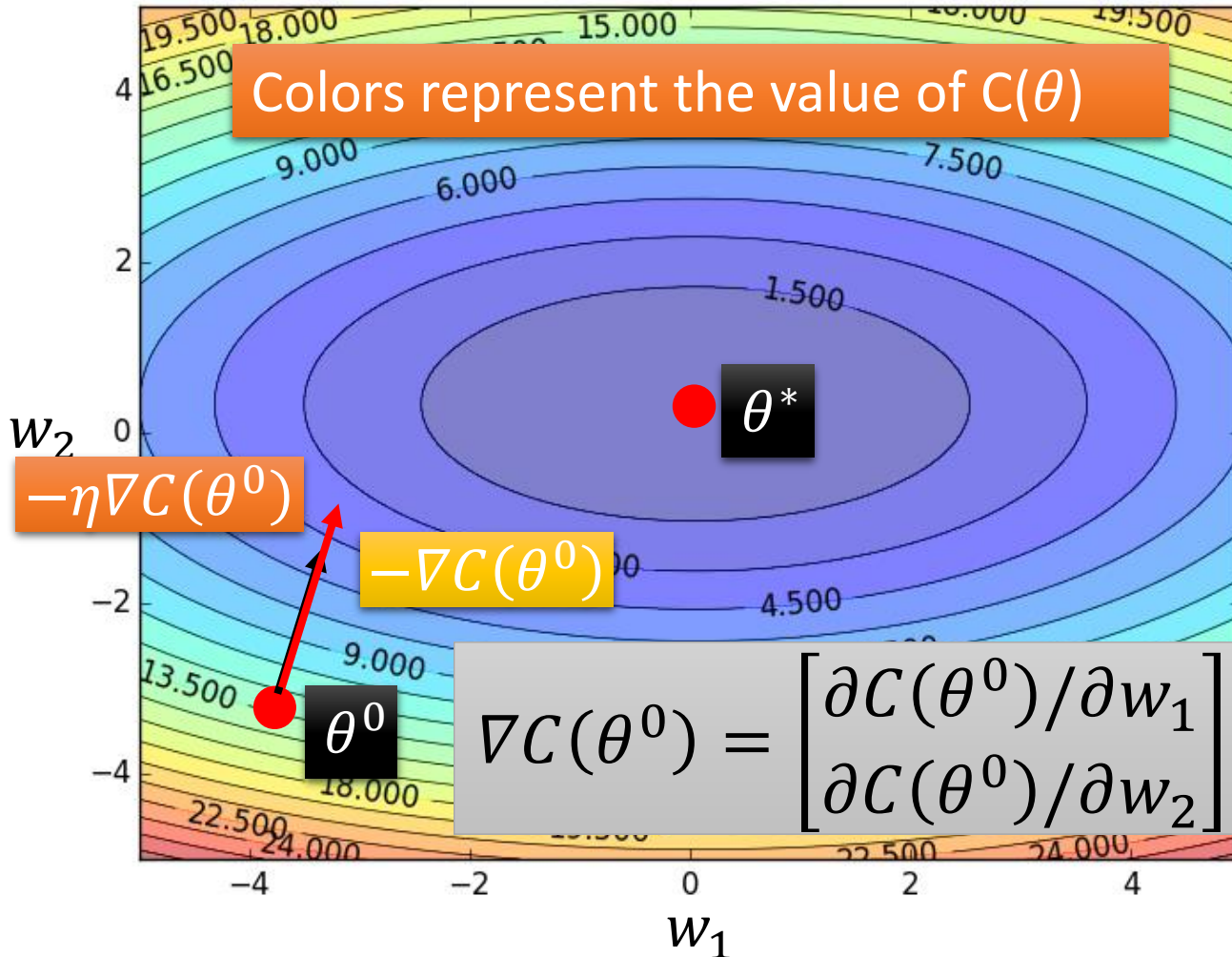
OPTIMIZATION

Gradient Descent

Assume there are only two parameters w_1 and w_2 in a network.

$$\theta = \{w_1, w_2\}$$

Error Surface



Randomly pick a starting point θ^0

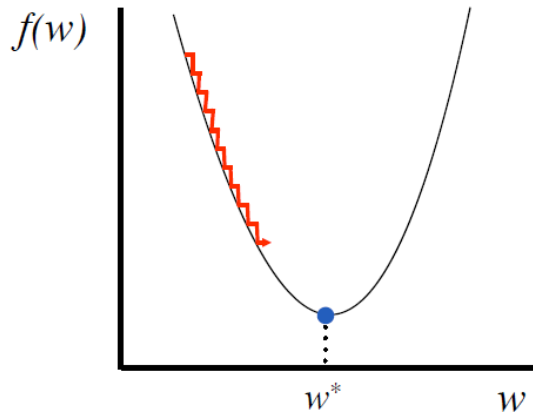
Compute the negative gradient at θ^0

➡ $-\nabla C(\theta^0)$

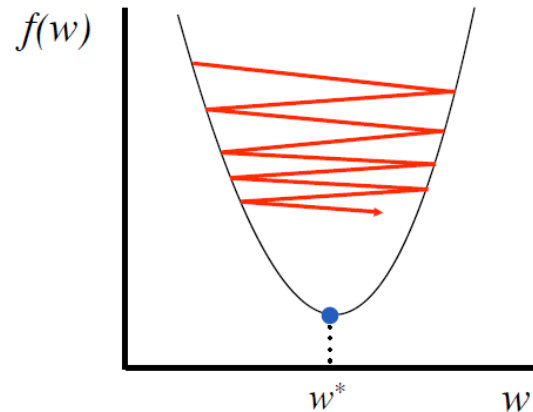
Times the learning rate η

➡ $-\eta \nabla C(\theta^0)$

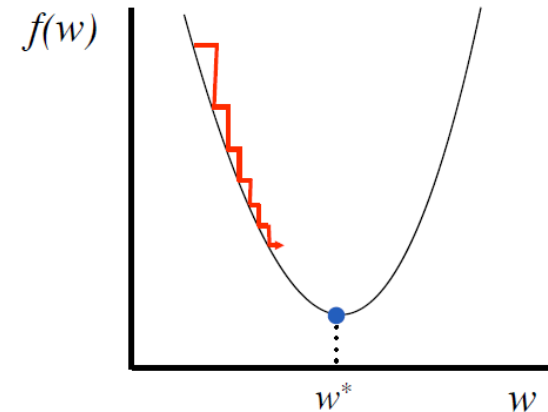
Choosing Step Size



Too small: converge very slowly



Too big: overshoot and even diverge



Reduce size over time

Theoretical convergence results for various step sizes

A common step size is $\alpha_i = \frac{\alpha}{n\sqrt{i}}$

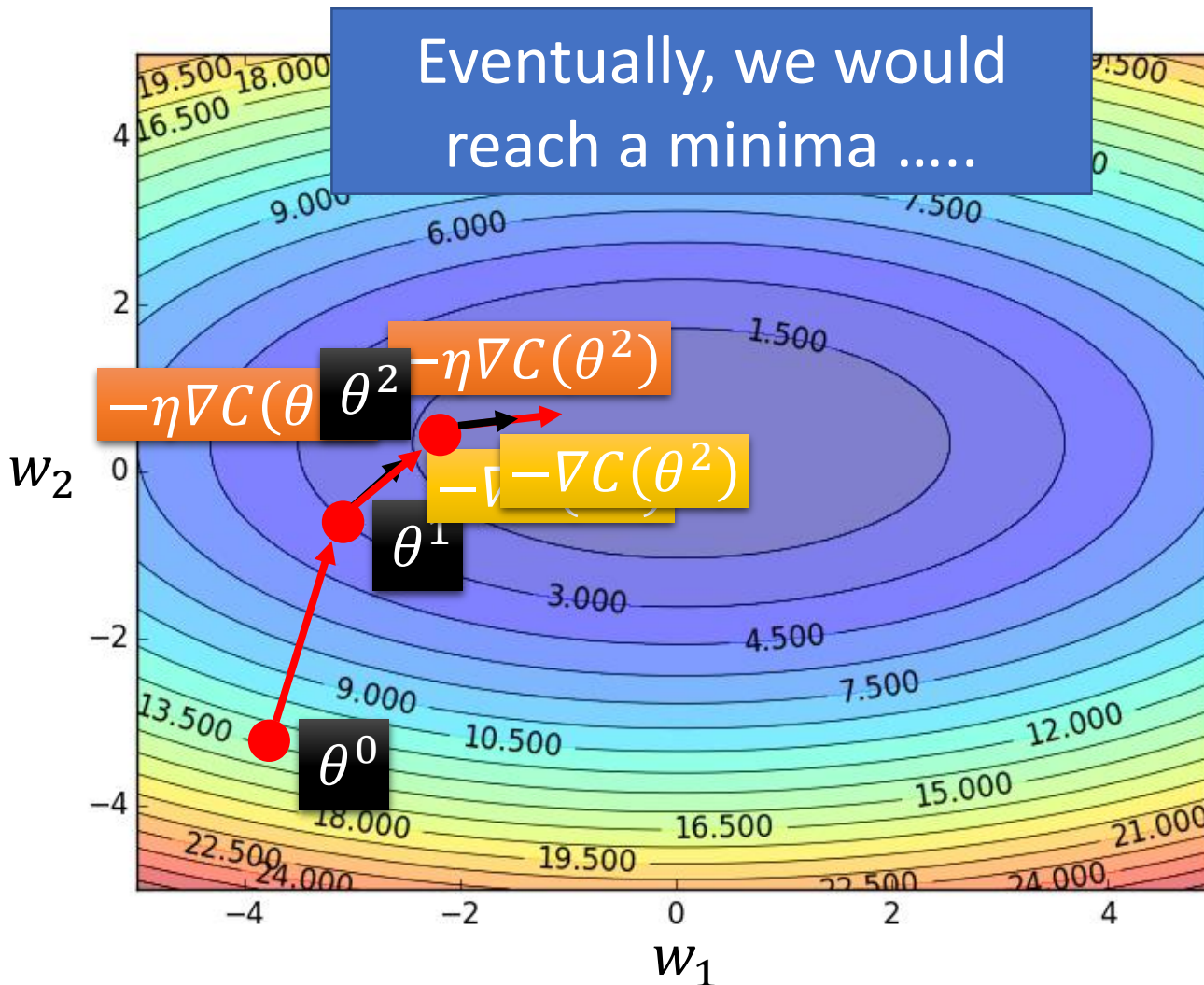
Training Points
 α — Constant
Iteration #

Source:



BerkeleyX: CS190.1x Scalable Machine Learning

Gradient Descent



Randomly pick a starting point θ^0

Compute the negative gradient at θ^0

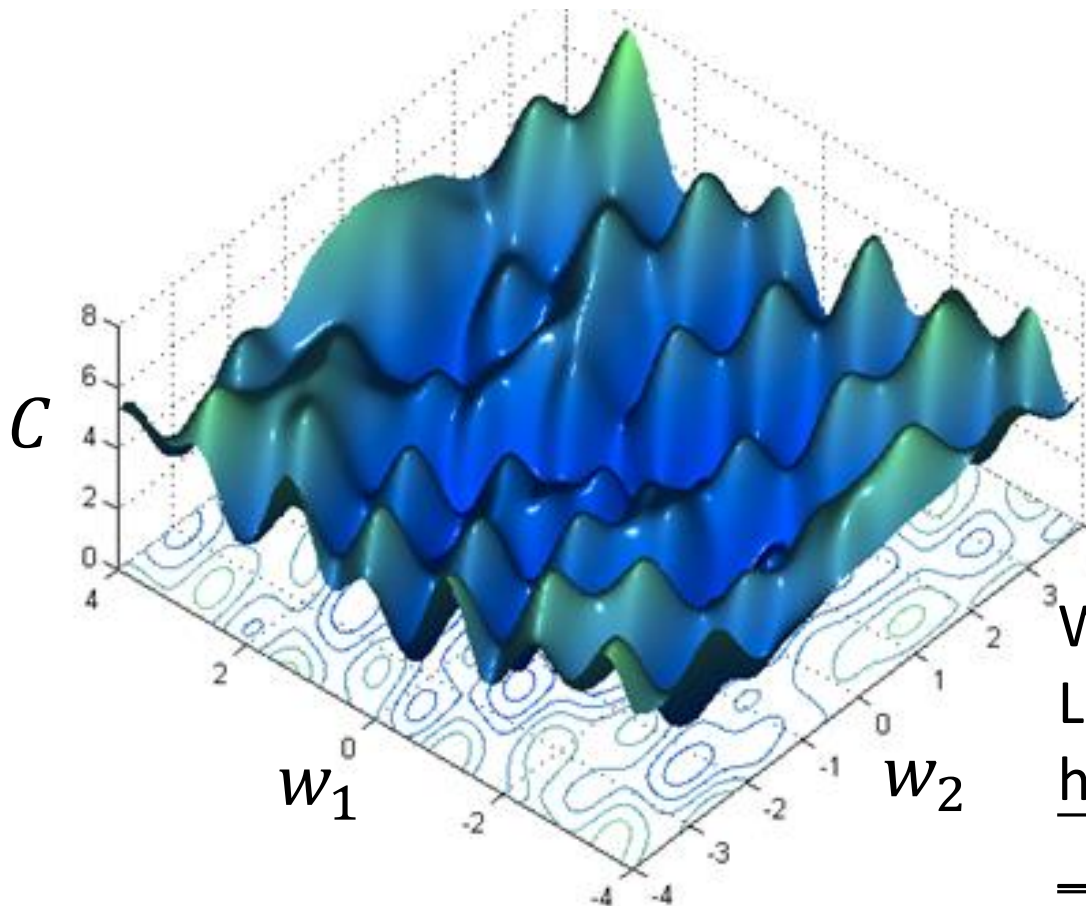
$\rightarrow -\nabla C(\theta^0)$

Times the learning rate η

$\rightarrow -\eta \nabla C(\theta^0)$

Local Minima

- Gradient descent never guarantee global minima



Different initial
point θ^0

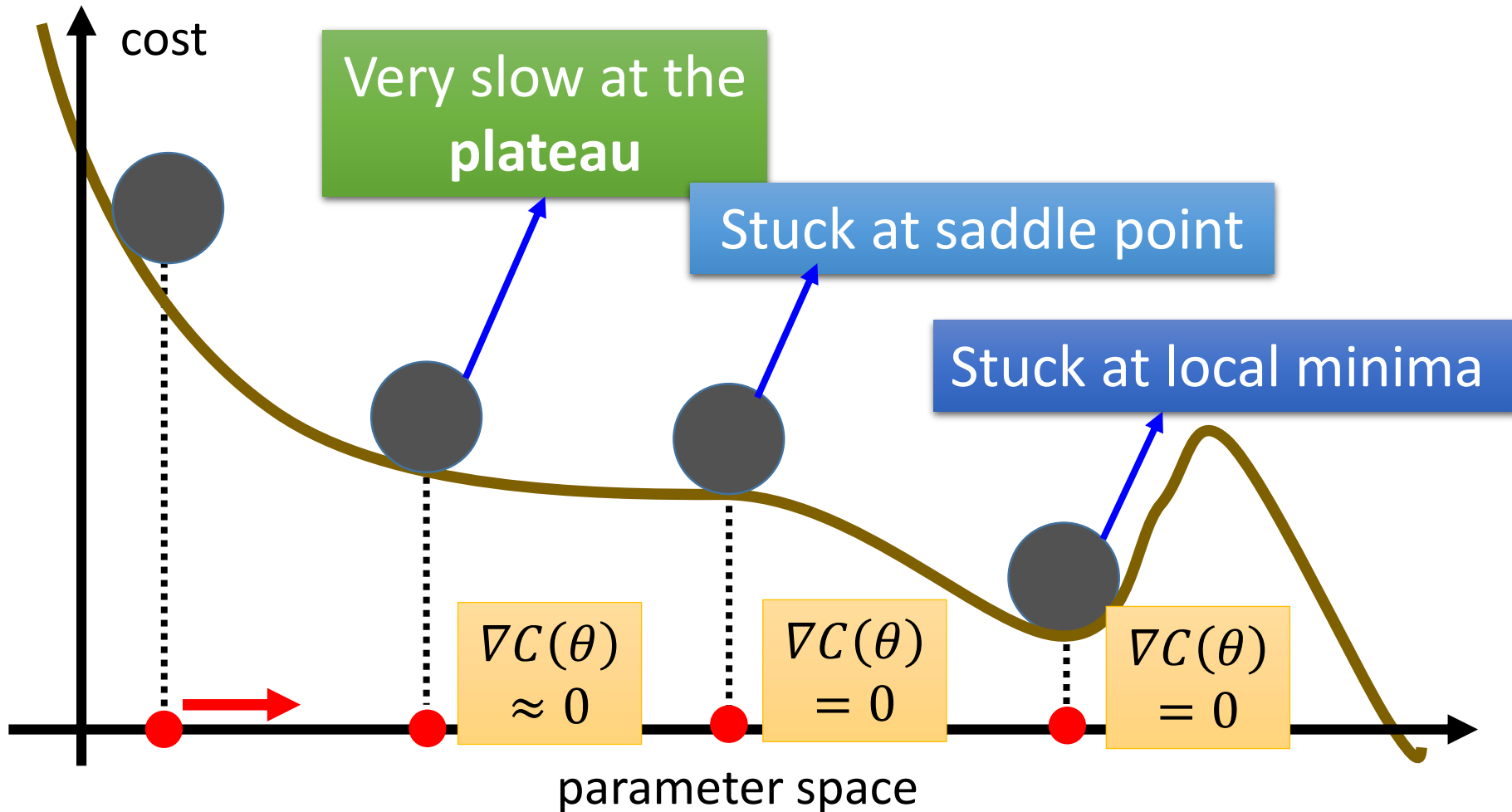


Reach different minima,
so different results

Who is Afraid of Non-Convex
Loss Functions?

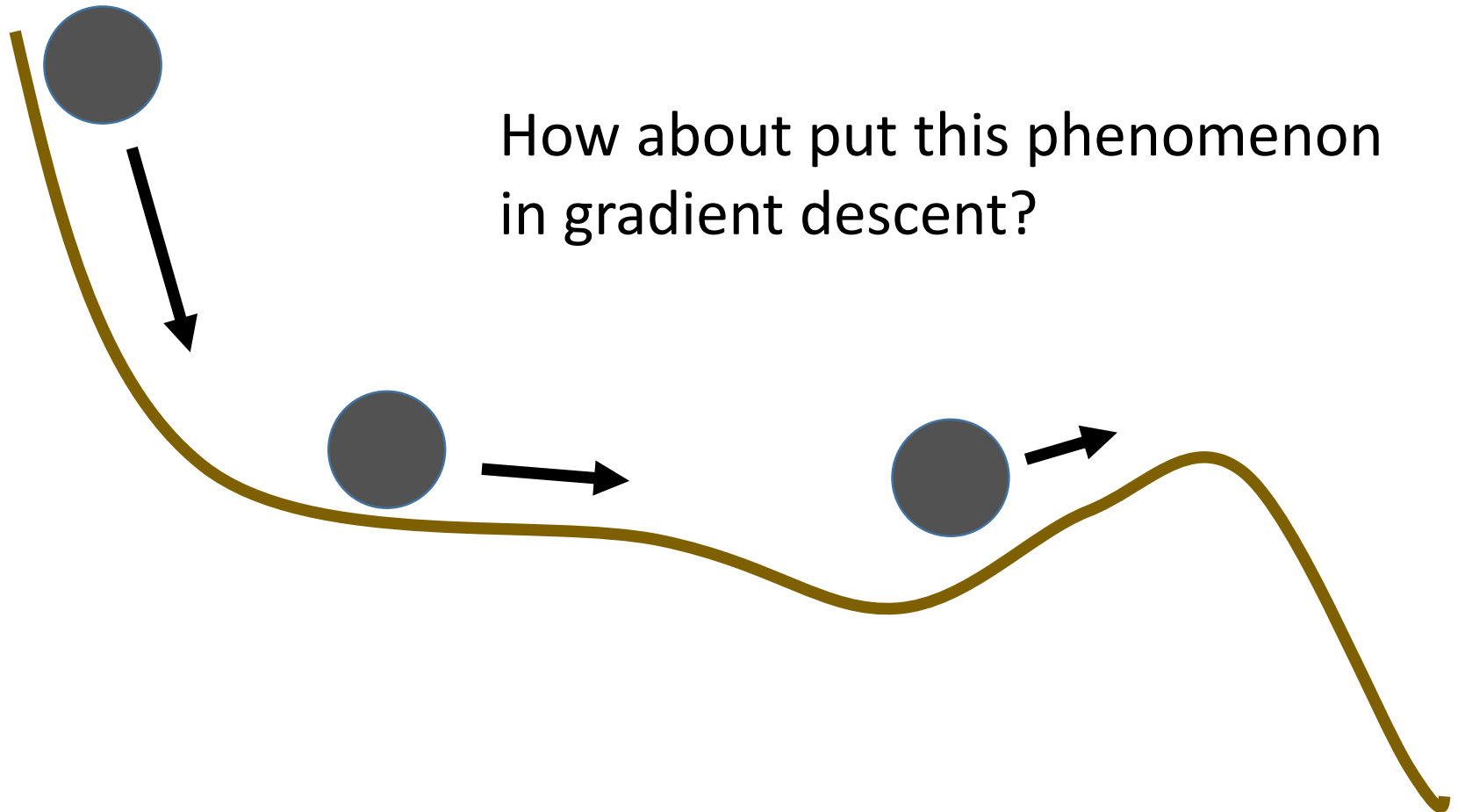
http://videolectures.net/eml07_lecun_wia/

Besides local minima



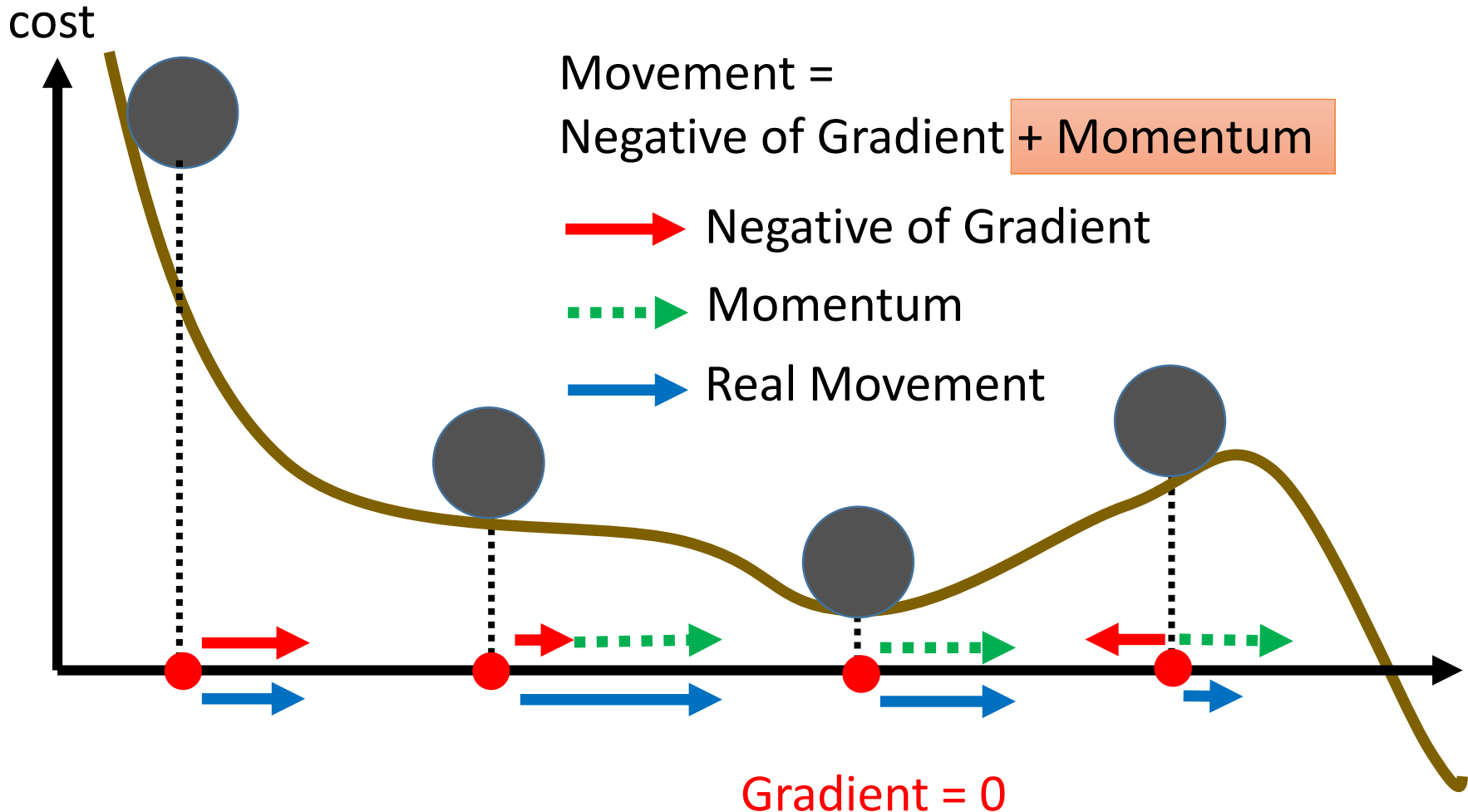
In physical world

- Momentum



Momentum

Still not guarantee reaching global minima, but give some hope

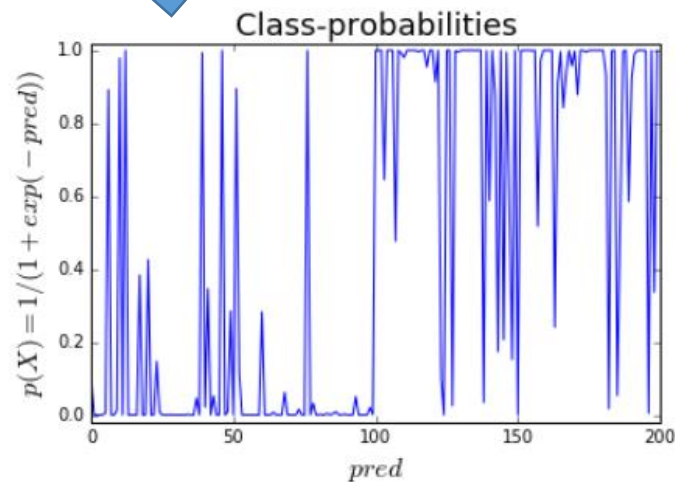
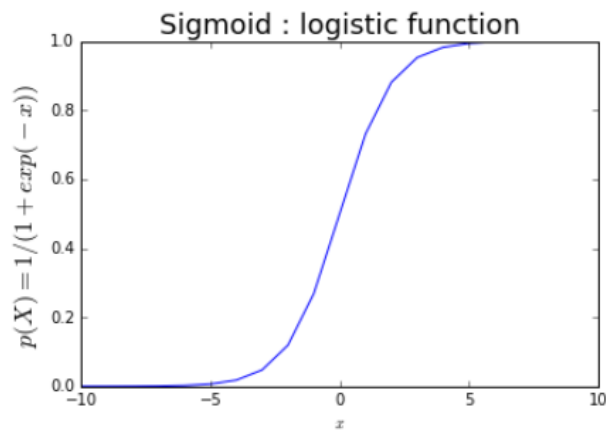
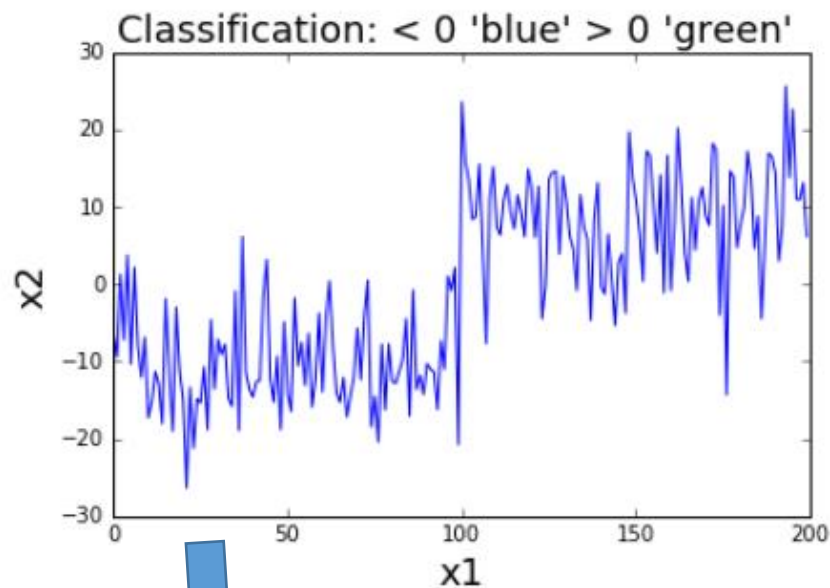
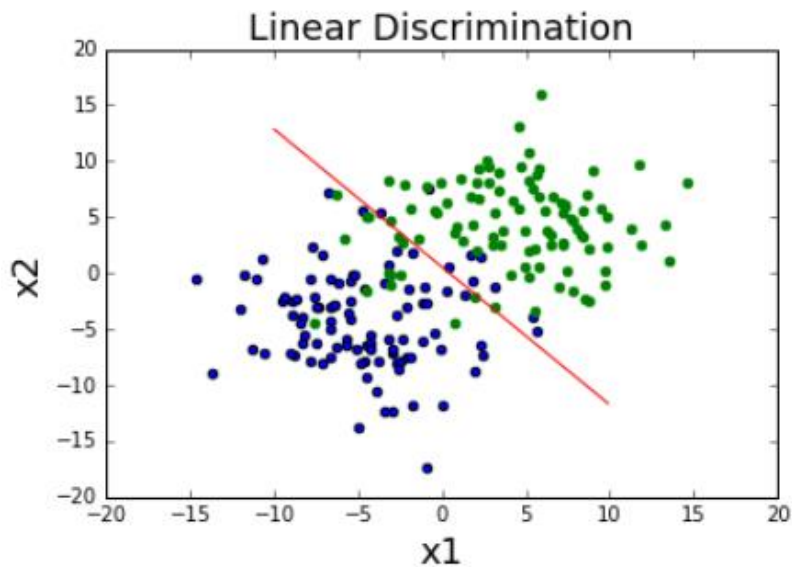


What Cost Function ?

- Classification error (number of errors)? How derivate? gradient?
- Better think on terms of **probability** & look for nice ways to derivate...

From scores to probabilities:

- Binary classification is easy: use sigmoid or logistic function



Cost function:

How to measure accuracy?

Cross-entropy and Maximum Likelihood Estimation

Intuitive approach for a two-class (binary) classifier

Take a data set:

$\{x_n, l_n\}$ of N vectors x_n belonging to two classes (labels $l_n = 0, 1$)

Consider our classifier as an estimator of the conditional probability:

- $p(l_n = 1|x_n) = p_{1n}$
- ...and so... $p(l_n = 0|x_n) = 1 - p_{1n}$

A PERFECT classifier

- If $l_n=1$ then $p(l_n = 1|x_n) = p_{1n}=1$
- If $l_n=0$ then $p(l_n = 1|x_n) = p_{1n}=0$

In a single expression:

$$p_{1n}^{l_n} (1 - p_{1n})^{(1-l_n)}$$

In a perfect classifier this must be:
a “probability” always 1 for every data $n!!!$

A PERFECT classifier

All the probabilities/likelihood for each sample should be 1

Maximum Likelihood estimation of θ

$$\theta_{ML} = \prod_{n=1}^N p_{1n}^{l_n} (1 - p_{1n})^{(1-l_n)}$$

From that a Cost Function can be obtained taking: **the negative logarithm = Cross-entropy**

$$\text{Cross entropy} = - \sum_{n=1}^N (l_n \log(p_{1n}) + (1 - l_n) \log(1 - p_{1n}))$$

Cross entropy loss function

From Maximum Likelihood to Minimum Cross-entropy

$$\text{Cross entropy} = - \sum_{n=1}^N (l_n \log(p_{1n}) + (1 - l_n) \log(1 - p_{1n}))$$

Shannon entropy $H(\cdot)$: of uncertainty in a probability distribution P

$$H(x) = -E_{x \sim P} [\log(P(x))]$$

Cross-entropy (two distributions : P_{data} and P_{model})

$$-E_{x \sim P_{data}} [\log(P_{model}(x))]$$

Cross entropy loss function

ANOTHER “PRACTICAL” POINT OF VIEW:

$$\text{Cross entropy} = - \sum_{n=1}^N (l_n \log(p_{1n}) + (1 - l_n) \log(1 - p_{1n}))$$

HEAVILY PENALIZES gross errors:

- $l_n = 1$ and p_{1n} close to 0 **$\log(0)!!$**
- $l_n = 0$ and p_{1n} close to 1 !! as $(1 - p_{1n})$ close to 0 **$\log(0)!!$**

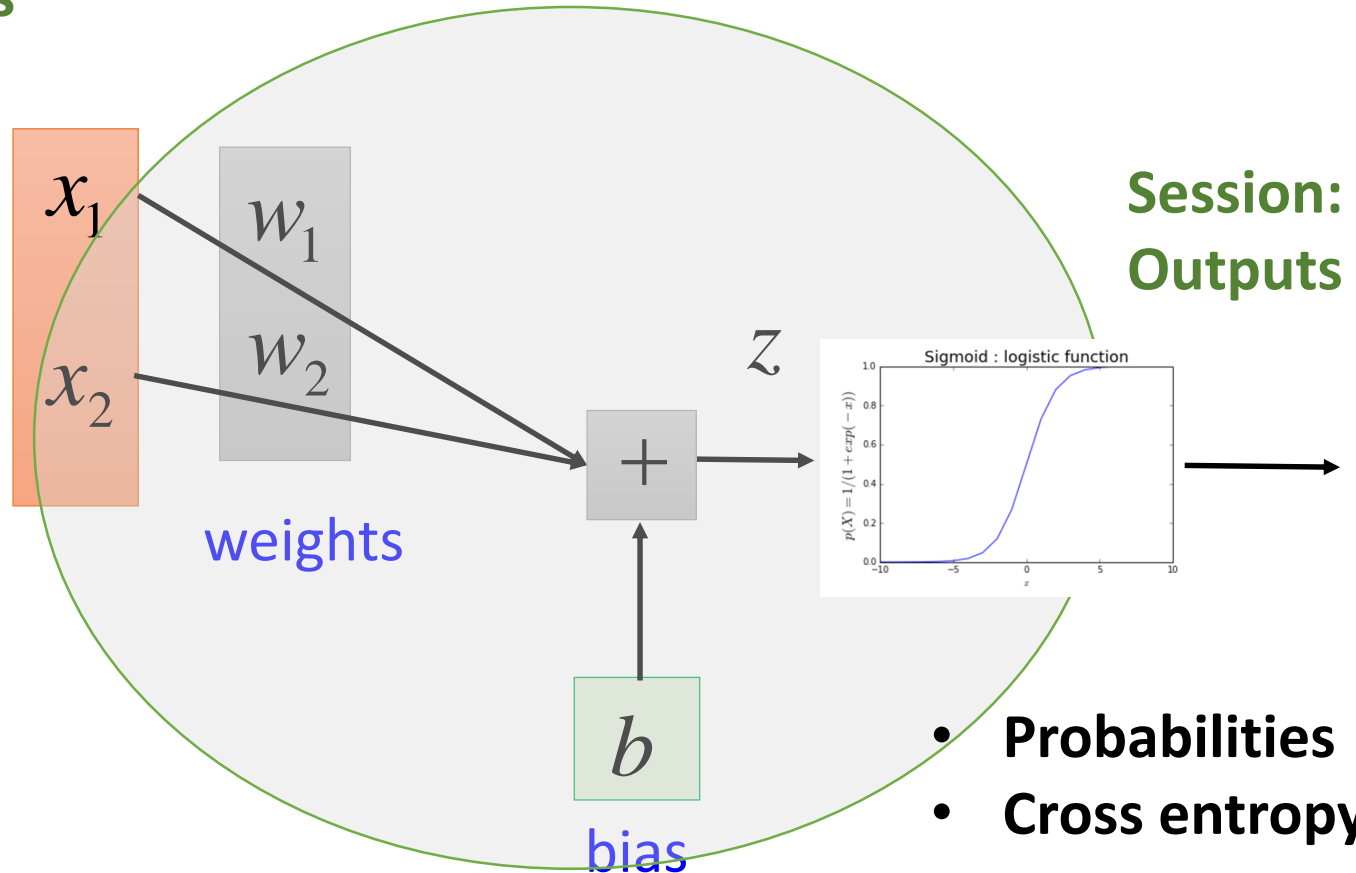
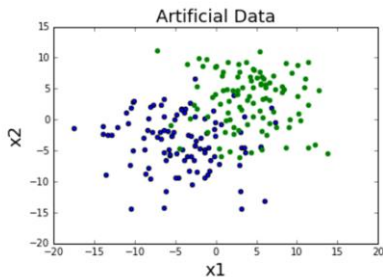


Now let's feed our Training data to a simple linear classifier

$$x_2 + 1.23 * x_1 - 0.55 = z$$

Inputs

Train data





```
W=tf.constant([[1.23], [1.0]],name="weights")  
b=tf.constant(-0.55,name="bias")
```

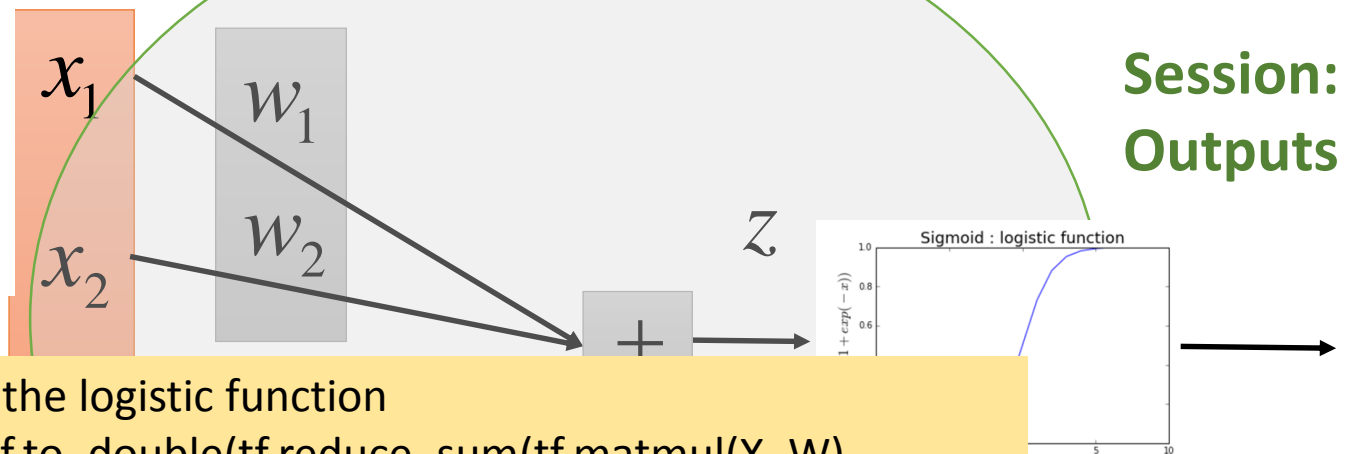
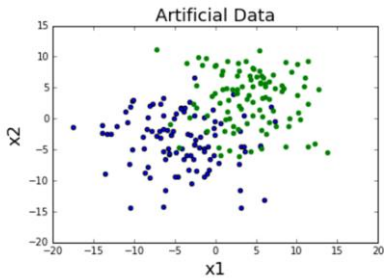
Inputs

$$x_2 + 1.23 * x_1 - 0.55 = z$$

```
X = tf.placeholder("float", shape=[None, 2])
```

Train data

```
labels = tf.placeholder("float", shape=[None])
```



Predictor is now the logistic function

```
pred = tf.sigmoid(tf.to_double(tf.reduce_sum(tf.matmul(X, W),  
axis=[1]) + b))
```

Cost function is cross-entropy

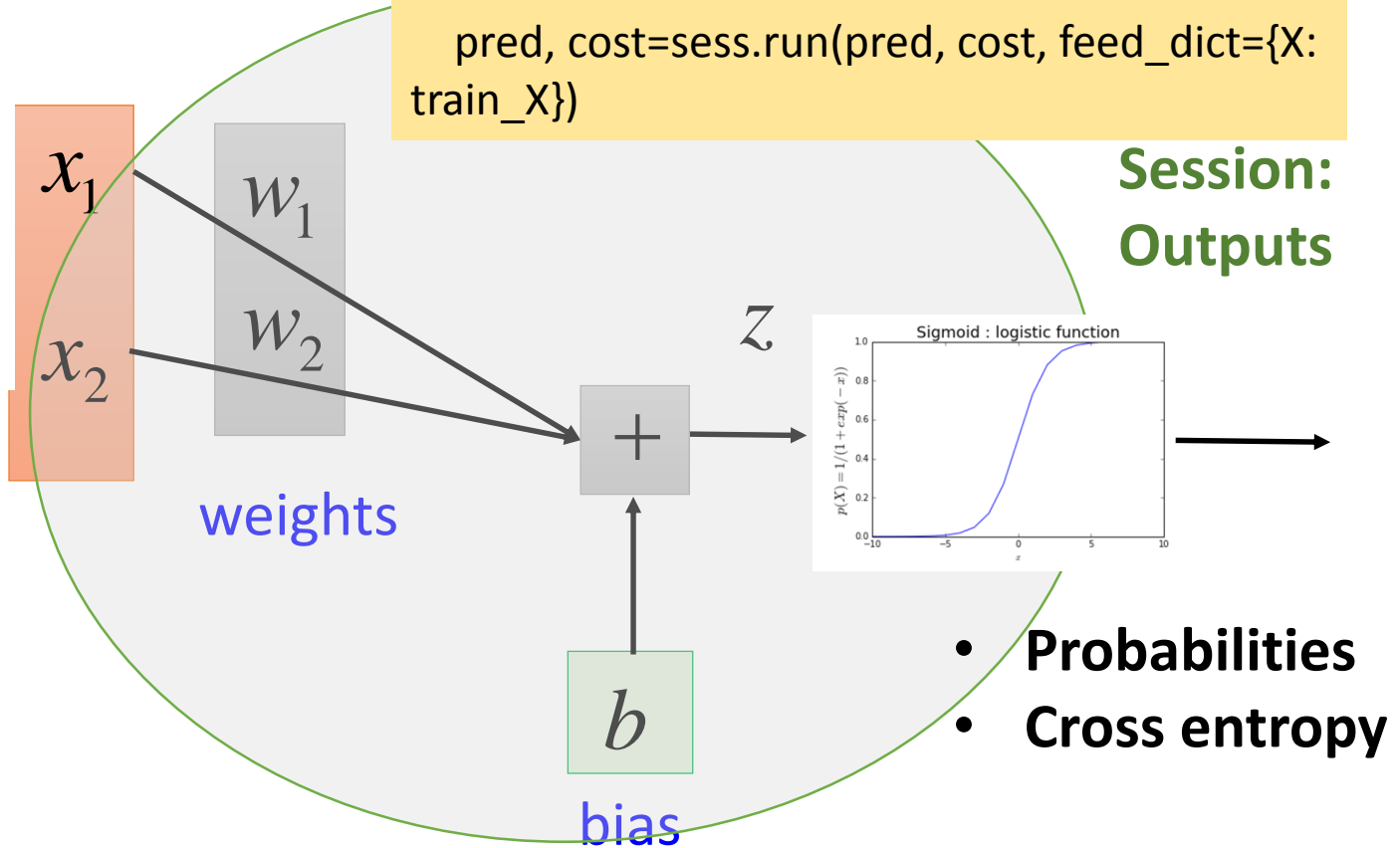
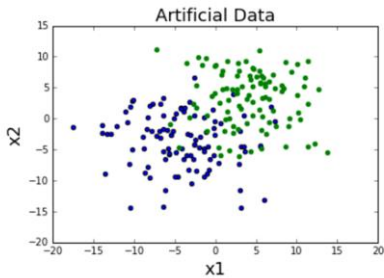
```
cost = -tf.reduce_sum(tf.to_double(labels) * tf.log(pred) + (1-  
tf.to_double(labels)) * tf.log(1-pred))
```

Probabilities
Cross entropy



Inputs

Train data



with tf.Session() as sess:

```
sess.run(init)
```

```
pred, cost=sess.run(pred, cost, feed_dict={X:  
train_X})
```



**Check that our results are the
same as before**

Now let's train!:



- **W** and **b** will be variables (tensors)

```
W = tf.Variable(tf.zeros([2, 1], "float"), name="weight")
```

```
b = tf.Variable(tf.zeros([1], "float"), name="bias")
```

- **and we use cross-entropy and gradient descend**

```
cost = -tf.reduce_sum(tf.to_double(labels) * tf.log(pred) + (1-  
tf.to_double(labels)) * tf.log(1-pred))
```

```
# Gradient descent
```

```
learning_rate = 0.001
```

```
optimizer = tf.train.GradientDescentOptimizer(learning_rate).minimize(cost)
```

- **We have used all the training data several runs (Epochs)**



```
with tf.Session() as sess:  
    sess.run(init)
```

```
# We can Run the optimization algorithm several times
```

```
for i in range(100):
```

```
    cost_out,W_out,b_out,_=sess.run([cost, W,b, optimizer], feed_dict={X:  
train_X, labels: train_labels})
```

```
    print("Epoch : %d Cost= %s"%(i,cost_out))
```

```
    print(W_out)
```

```
    print(b_out)
```


When large amounts of data divide data into mini-batches



- Most optimization algorithms converge much faster (in terms of total computation, not in terms of number of updates) if they are allowed to rapidly compute approximate estimates of the gradient rather than slowly computing the exact gradient.
- Another consideration motivating statistical estimation of the gradient from a small number of samples is redundancy in the training set.



- Optimization algorithms that use the entire training set are called **batch or deterministic gradient**
- **Stochastic Gradient Descend (SGD):** Optimization algorithms that use only a single example at a time
- Most algorithms used for deep learning fall somewhere in between: **minibatch or minibatch stochastic methods**



Confusing terminology:

- The word “batch” is also often used to describe the minibatch used by minibatch stochastic gradient descent.
- It is very common to use the term “batch size” to describe the size of a minibatch

See more details on how choosing minibatch size on Deep Learning Book
(Chap 8 : Optimization)

.... See details in Notebook
... and practice with it....

- Random initialization of variables
- Stepsize
- Optimizers
- Interactive Session
- Tf Debugging?

Deep Learning Seminar (materials)

Deep Learning using TensorFlow and TensorFlow-Slim

Dipendra Jha Northwestern University

dipendra009@gmail.com <https://www.linkedin.com/in/dipendra009>

Deep Learning courses

Prof. Hung-yi Lee National Taiwan University (NTU) Taipei

Introduction to Deep Learning

[Yingyu Liang](#) Princeton University

<http://jrmeyer.github.io/tutorial/2016/02/01/TensorFlow-Tutorial.html>

<http://www.psi.toronto.edu/~jimmy/ece521/Tut1.pdf>