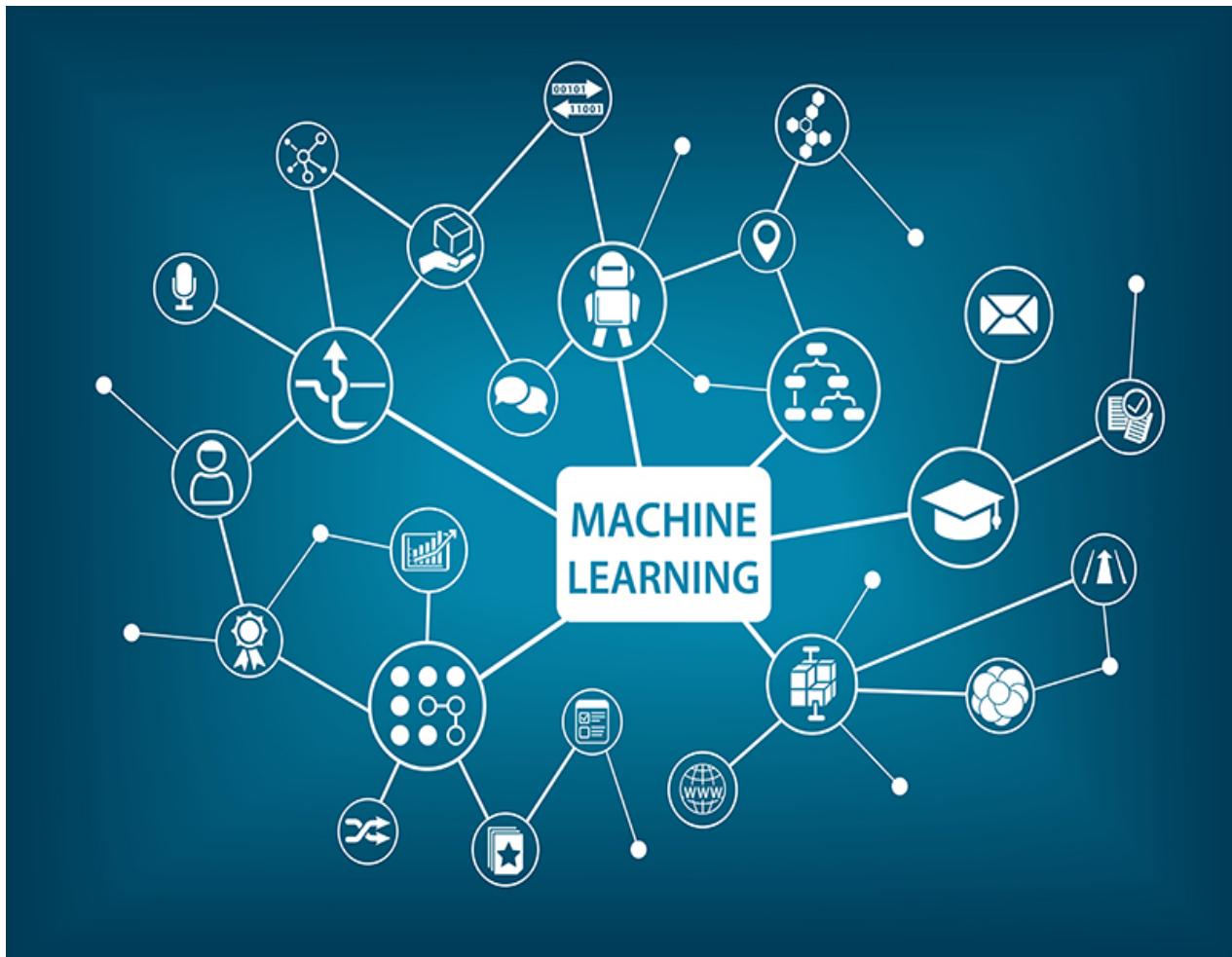


Introduction to machine learning and regression



András Horváth

This presentation is available at:

[http://users.itk.ppke.hu/~horan
/big_data](http://users.itk.ppke.hu/~horan/big_data)

The code is available at:

[https://databricks-prod-cloudfront.cloud.databricks.com/
public/4027ec902e239c93eaaa8714f173bcfc/2217267761988785
/1705666210822011/2819141955166097/latest.html](https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/2217267761988785/1705666210822011/2819141955166097/latest.html)

Machine learning, machine lintelligence

What is intelligence?



Intelligent life can't be all that common because it's really rare on Earth and especially since we define ourselves to be intelligent. But in the eyes of an alien coming here who has the technology to make it here, they might observe us and conclude that there's no sign of intelligent life on Earth.

— Neil deGrasse Tyson —

AZ QUOTES

SOMETIMES I THINK
THE SUREST SIGN THAT
INTELLIGENT LIFE
EXISTS ELSEWHERE IN
THE UNIVERSE IS THAT
NONE OF IT HAS TRIED
TO CONTACT US.



Machine learning, machine lintelligence

What is intelligence?

The ability to acquire and apply knowledge and skills.



Machine learning, machine lintelligence

What is intelligence?

The ability to acquire and apply knowledge and skills.

**1N73LL1G3NC3
15 7H3
4B1L17Y
70 4D4P7 70
CH4NG3.
-573PH3N H4WK1NG**

Machine learning, machine intelligence

What is intelligence?

The ability to acquire and apply knowledge and skills

Intelligence is the ability to adapt to change

Providing computers the ability to learn without being explicitly programmed:

Involves: programming, Computational statistics, mathematical optimization, image processing, natural language processing etc...



Machine learning, machine intelligence

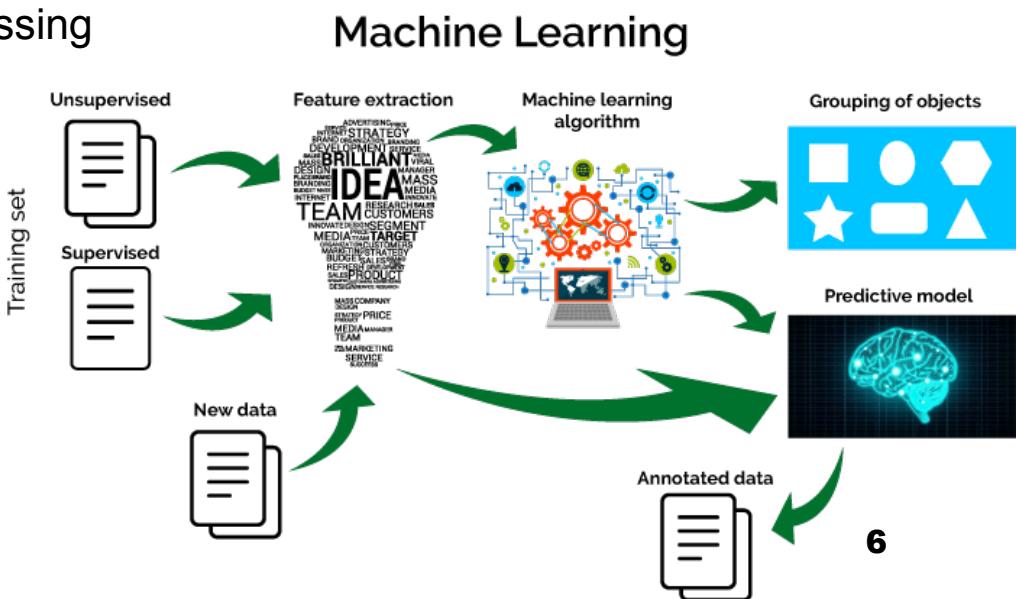
What is intelligence?

The ability to acquire and apply knowledge and skills

Intelligence is the ability to adapt to change

Providing computers the ability to learn without being explicitly programmed:

Involves: programming, Computational statistics, mathematical optimization, image processing, natural language processing etc...



Conquests of machine learning

1952 Arthur Samuel (IBM): First machine learning program playing checkers

Arthur Samuel coined the term „machine learning”



Conquests of machine learning

1952 Arthur Samuel (IBM): First machine learning program playing checkers

1997 IBM Deep Blue Beats Kasparov:

First match (1996 Nov): Kasparov–Deep Blue (4–2)

Second Match (1997 May): Deep Blue–Kasparov (3½–2½)



Conquests of machine learning

1952 Arthur Samuel (IBM): First machine learning program playing checkers

1997 IBM Deep Blue Beats Kasparov:

2011 IBM Watson: Beating human champions in Jeopardy

It's a 4-letter term for a summit; the first 3 letters mean a type of simian : Apex

4-letter word for a vantage point or a belief : View

Music fans wax rhapsodic about this Hungarian's "Transcendental Etudes" : Franz Liszt

While Maltese borrows many words from Italian, it developed from a dialect of this semitic language : Arabic



Conquests of machine learning

1952 Arthur Samuel (IBM): First machine learning program playing checkers

1997 IBM Deep Blue Beats Kasparov:

2011 IBM Watson: Beating human champions in Jeopardy

2014 Deep face algorithm Facebook

Reached 97.35% accuracy

Human performance is around 97%

Who's in These Photos?

The photos you uploaded were grouped automatically so you can quickly label and notify friends in these pictures. (Friends can always untag themselves.)



Conquests of machine learning

1952 Arthur Samuel (IBM): First machine learning program playing checkers

1997 IBM Deep Blue Beats Kasparov:

2011 IBM Watson: Beating human champions in Jeopardy

2014 Deep face algorithm Facebook

2016 Alpha go: deep learning

Fan Hui (5-0)

Lee Sedol (4-1)

99.8% win rate against other Go programs



Types of machine learning

Supervised learning and unsupervised learning

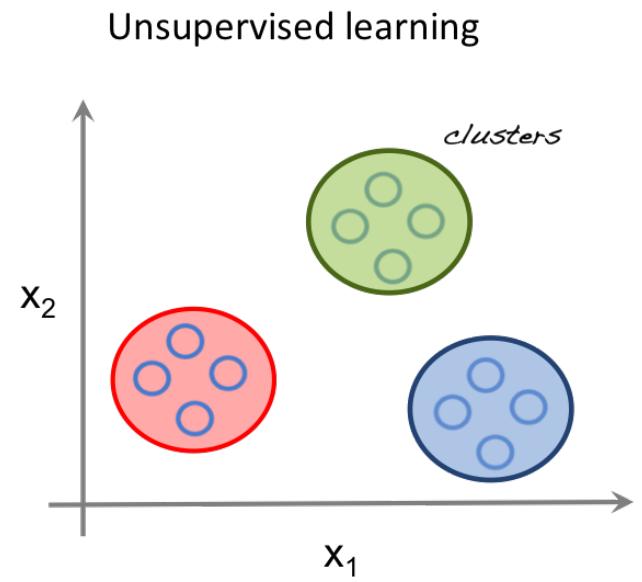
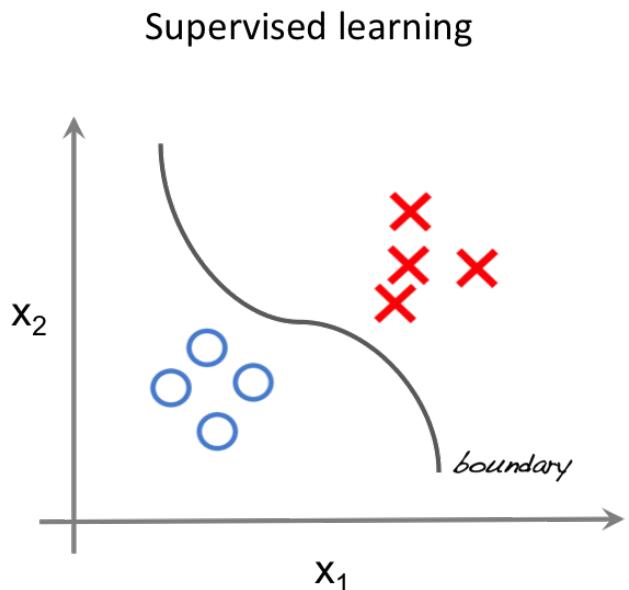
Types of machine learning

Supervised learning and unsupervised learning

Supervised:

We know the expected output we know the perfect, desired output on the training set

Classification and regression



Unsupervised:

All we have is data and no labels. We have to identify rules in the structure of the data

Clustering and association

Types of machine learning

Supervised learning and unsupervised learning

Supervised:

Classification and regression

Classification: A classification problem is when the output variable is a category, such as “red” or “blue” or “disease” and “no disease”.

Regression: A regression problem is when the output variable is a real value, such as “dollars” or “weight”.

Unsupervised:

Clustering and association

Clustering: A clustering problem is where you want to discover the inherent groupings in the data, such as grouping customers by purchasing behavior.

Association: An association rule learning problem is where you want to discover rules that describe large portions of your data, such as people that buy X also tend to buy Y.

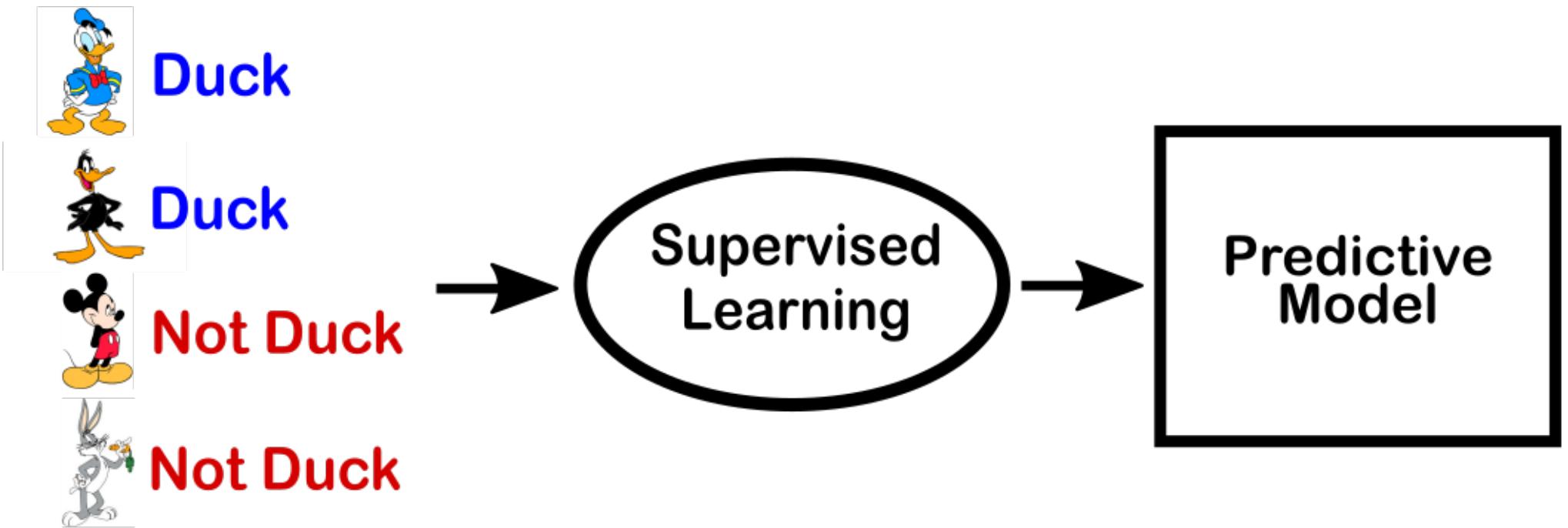
Semi-supervised: Some data is labeled but most of it is unlabeled and a mixture of supervised and unsupervised techniques can be used.

Supervised Classification

The most important conquests of deep learning comes from supervised classification.

There are three major improvements in this:

- New methods and technologies
- Powerful hardware for training (GPGPUs)
- Vast amount of data is available



Supervised Regression

Our aim is to predict a target value (Y) based on our data (X)

We have to find a model that explains the rules how Y can be derived from X $m(X)=Y$

The perfect fit is usually not possible, because our world is not perfect:

We have a model which has flaws

There is noise in our observation

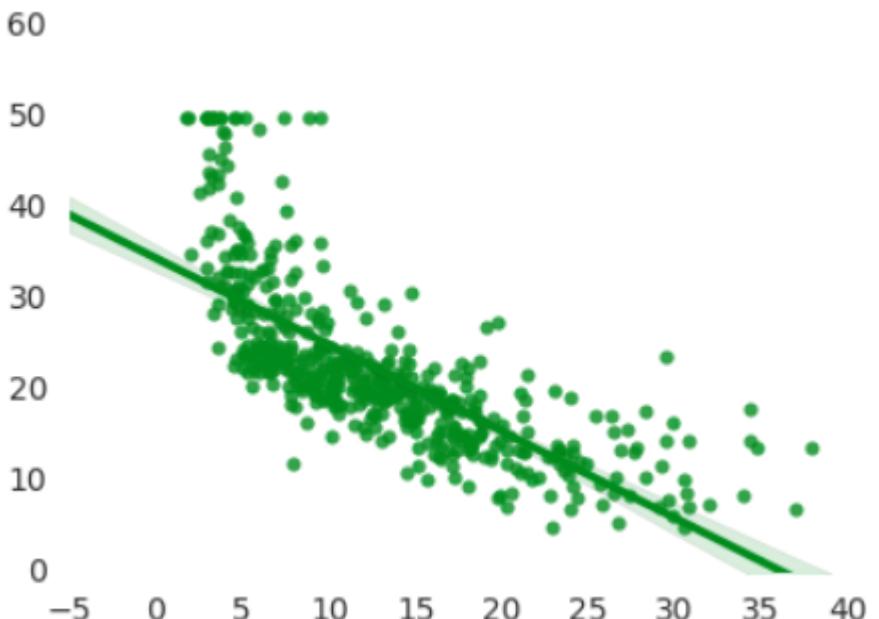
We have to fit a model which minimizes the error on the dataset:

L1 error

$$\sum_i |m(X_i) - Y_i|$$

L2 error

$$\sum_i (m(X_i) - Y_i)^2$$

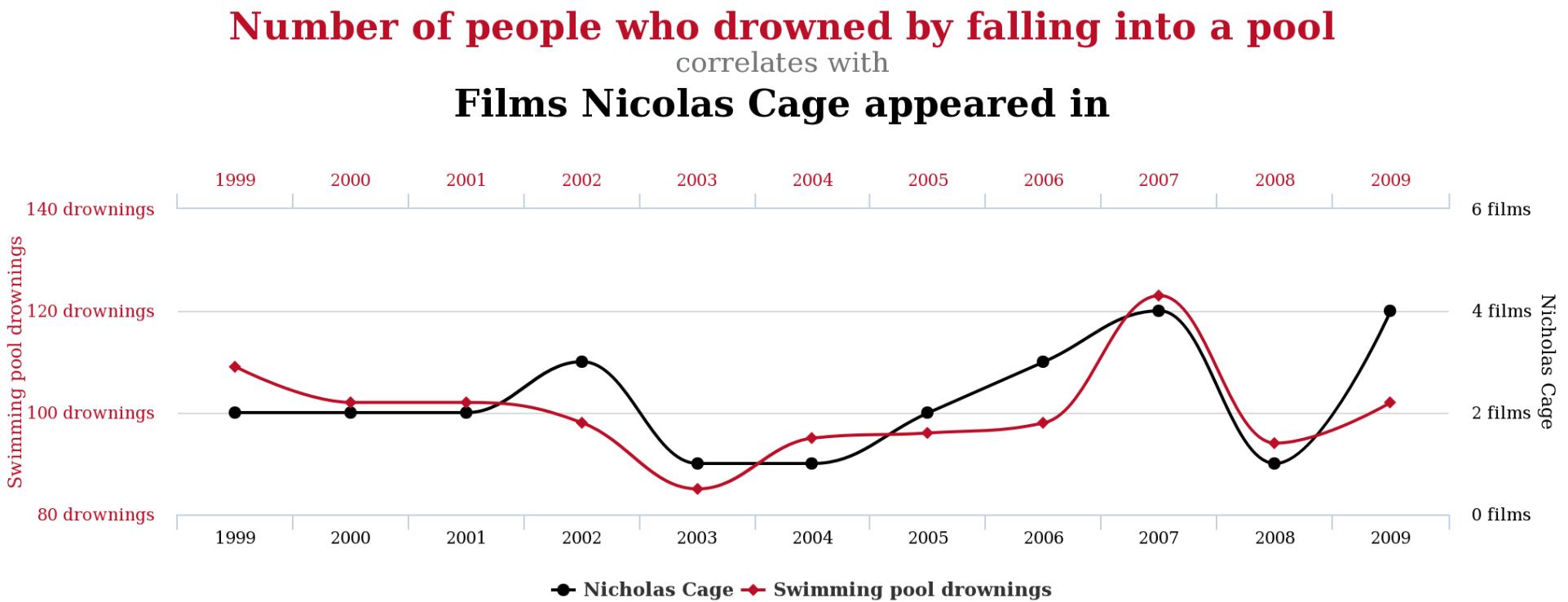


Correlation and cause effect relationship

Large correlation between the data does not necessary mean that one is caused by the other

One has to gather the large amount of data to ensure the coincidence by luck

One has to understand the data



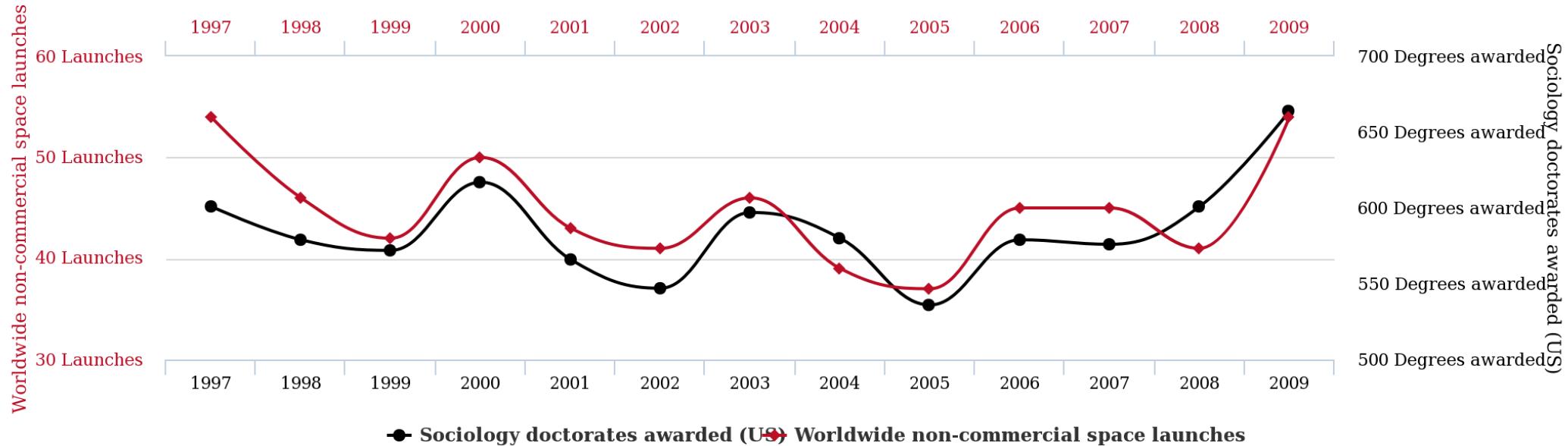
Correlation and cause effect relationship

Large correlation between the data does not necessary mean that one is caused by the other

One has to gather the large amount of data to ensure the coincidence by luck

One has to understand the data

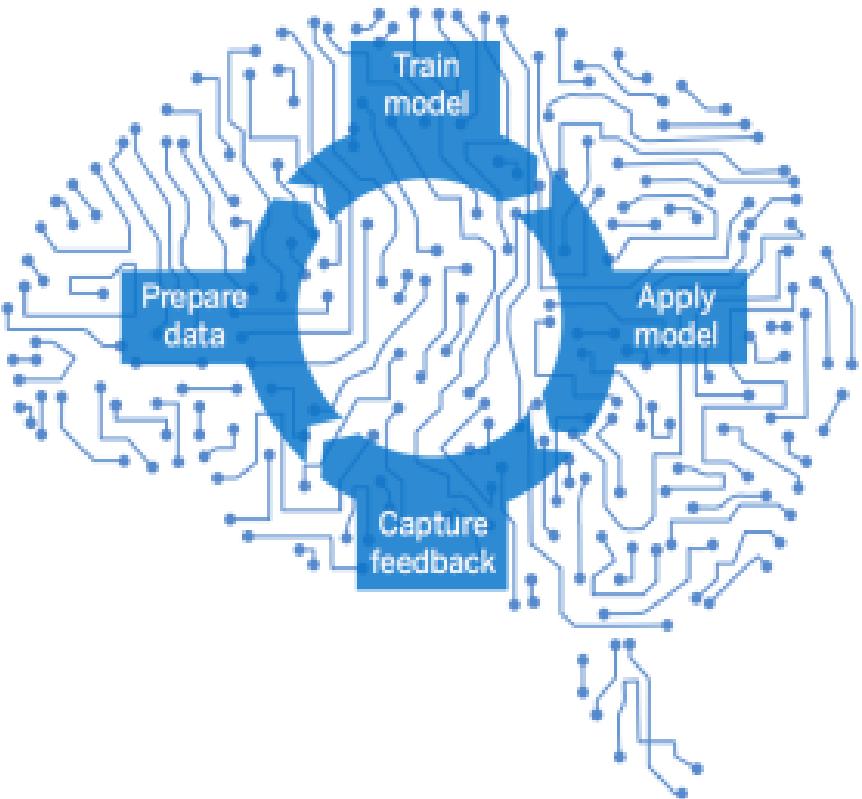
Worldwide non-commercial space launches
correlates with
Sociology doctorates awarded (US)



Steps of machine learning

Supervised learning and unsupervised learning:

- 1, Gather Data
- 2, Understand your data
- 3, Data preparation
- 4, Choosing a model
- 5, Training and evaluation
- 6, Parameter tuning
- 7, Go back to point 1



Gathreing data

Lets have an example dataset:

Boston house prices dataset:

A typical regression problem, good for trying out different machine learning models

There are many available solutions on this dataset

Standard dataset for model evaluation contained by Python libraries

Easy to load the data:

```
1 #gather data
2 # we will import modules which already contain the dataset
3 import sklearn
4 from sklearn.datasets import load_boston
5 #load the database
6 boston = load_boston()
```



Understanding the data

```
1 #understand the data
2 print(boston.DESCR) #contains the description of the dataset
3 print(boston.data.shape) # the shape of the dataset (rows and columns)
4 print(boston.keys())
5 X, Y = boston.data, boston.target
6 import pandas as pd
7 boston_df = pd.DataFrame(boston.data)
8 boston_df.columns = boston.feature_names
9 boston_df.describe()
```

The dataset contains 506 rows (different data points) and 14 columns

13 features describing the properties (X):

per capita crime rate by town, proportion of residential land zoned for lots over 25,000 sq.ft., proportion of non-retail business acres per town, Charles River dummy variable (= 1 if tract bounds river; 0 otherwise), nitric oxides concentration (parts per 10 million), average number of rooms per dwelling, proportion of owner-occupied units built prior to 1940, weighted distances to five Boston employment centers, index of accessibility to radial highways, full-value property-tax rate per \$10,000, pupil-teacher ratio by town, $1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town, % lower status of the population

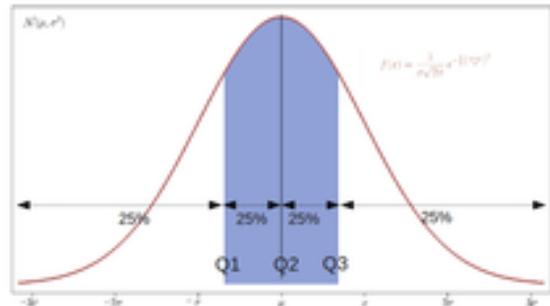
14th column: the target value (Y)

Our aim is to predict Y from X

Data preparation

Write out some values:

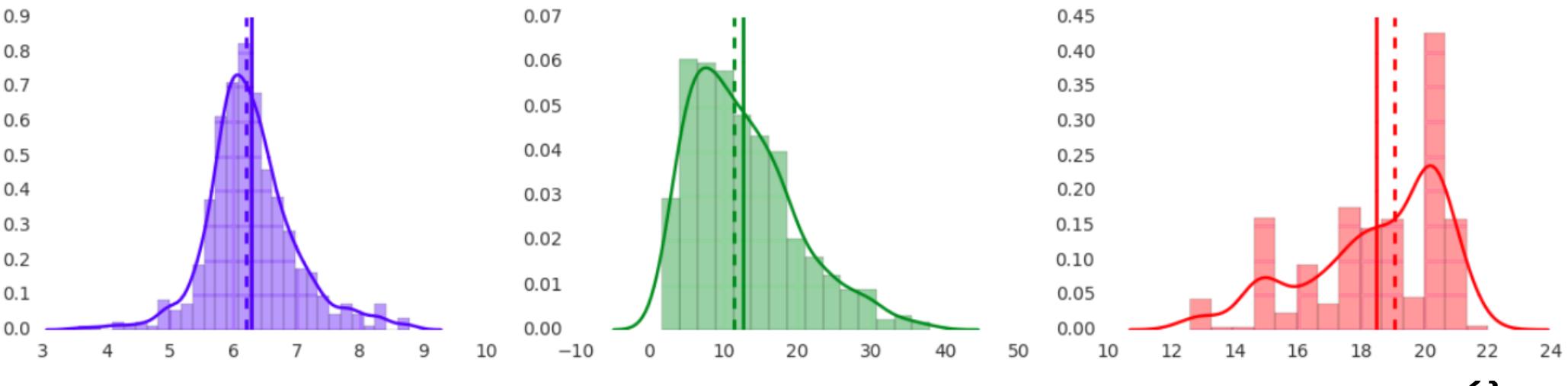
```
1 #check values
2 boston_df.min() #minimum values of the features
3 boston_df.max() #maximum values of the features
4 boston_df.median()#median values of the features
5 boston_df.quantile(0.5) #50 percent quantile of the features
```



Data preparation

Plot some values as histograms:

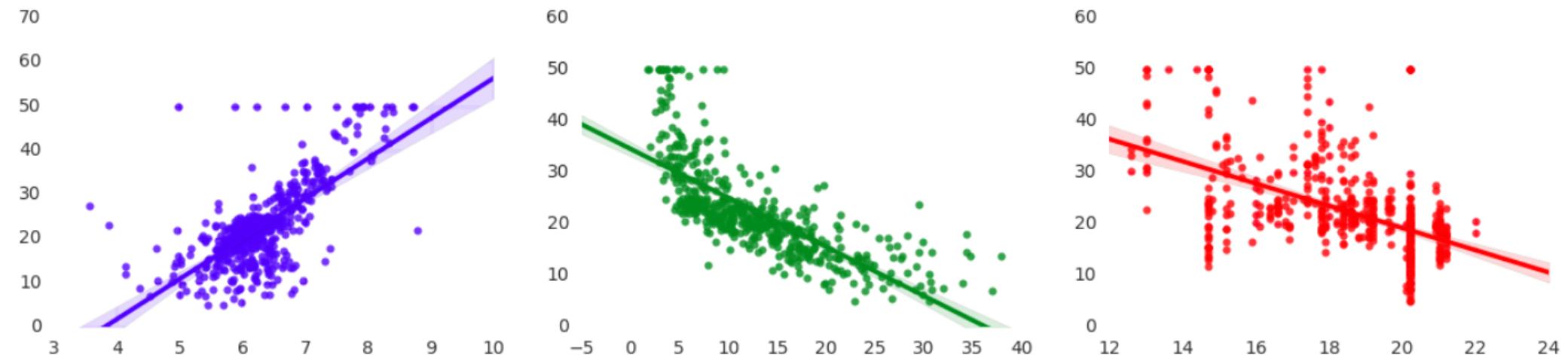
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 clr = ['blue', 'green', 'red']
4 fig, ax = plt.subplots(ncols=3, figsize=(15,3))
5
6 for i, var in enumerate(['RM', 'LSTAT', 'PTRATIO']):
7     sns.distplot(boston_df[var], color = clr[i], ax=ax[i])
8     ax[i].axvline(boston_df[var].mean(), color=clr[i], linestyle='solid', linewidth=2)
9     ax[i].axvline(boston_df[var].median(), color=clr[i], linestyle='dashed', linewidth=2)
10
11 display(fig)
```



Data preparation

Plot some values related to the average price:

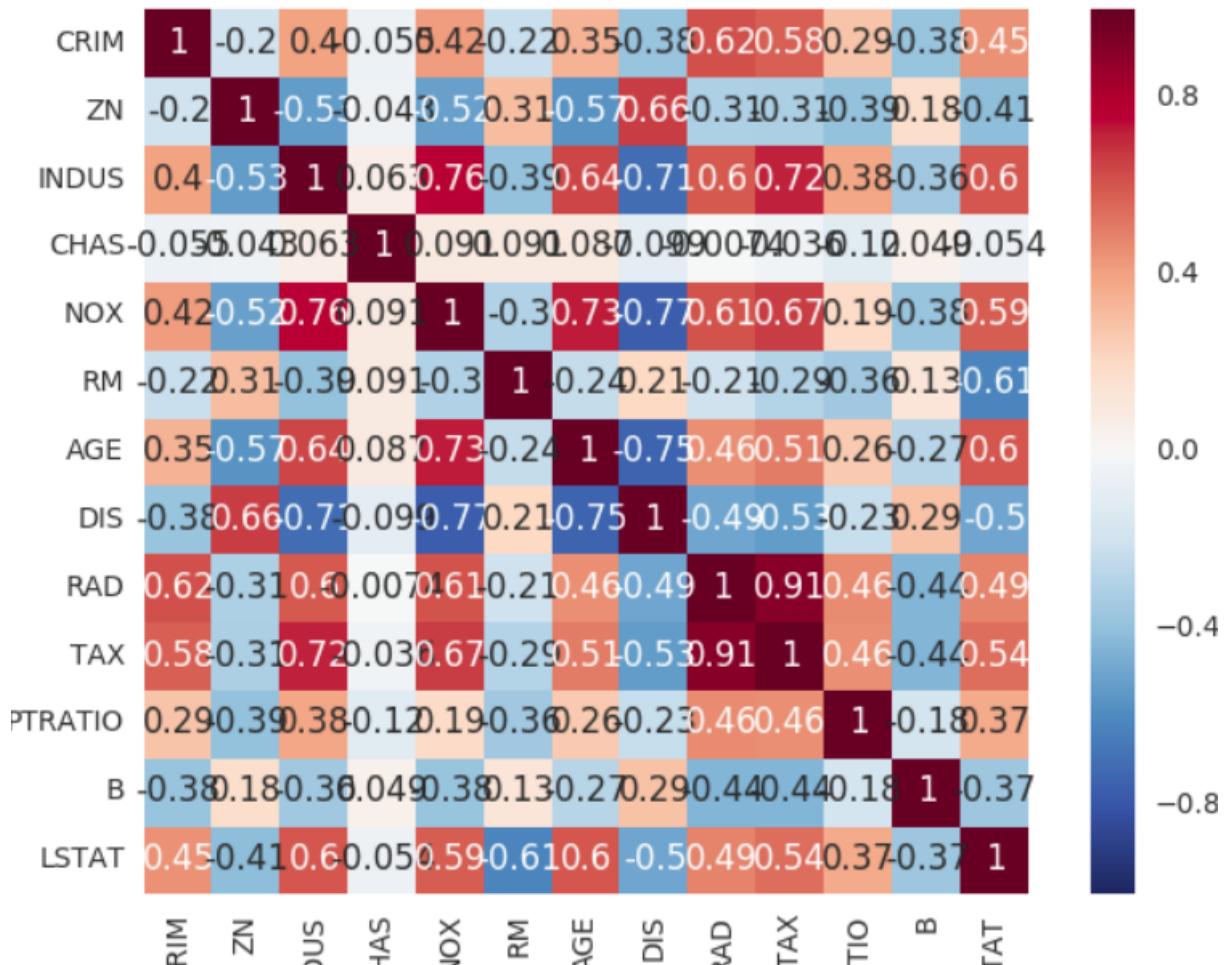
```
1 fig, ax = plt.subplots(ncols=3, figsize=(15,3))
2
3 for i, var in enumerate(['RM', 'LSTAT', 'PTRATIO']):
4     sns.regress(boston_df[var], boston.target, color=clr[i], ax=ax[i])
5     ax[i].set(ylim=(0, None))
6 display(fig)
```



Data preparation

Plot correlation between the different features:

```
1 fig,ax= plt.subplots()  
2 sns.heatmap(boston_df.corr(), square=True, annot=True)  
3 display(fig)
```



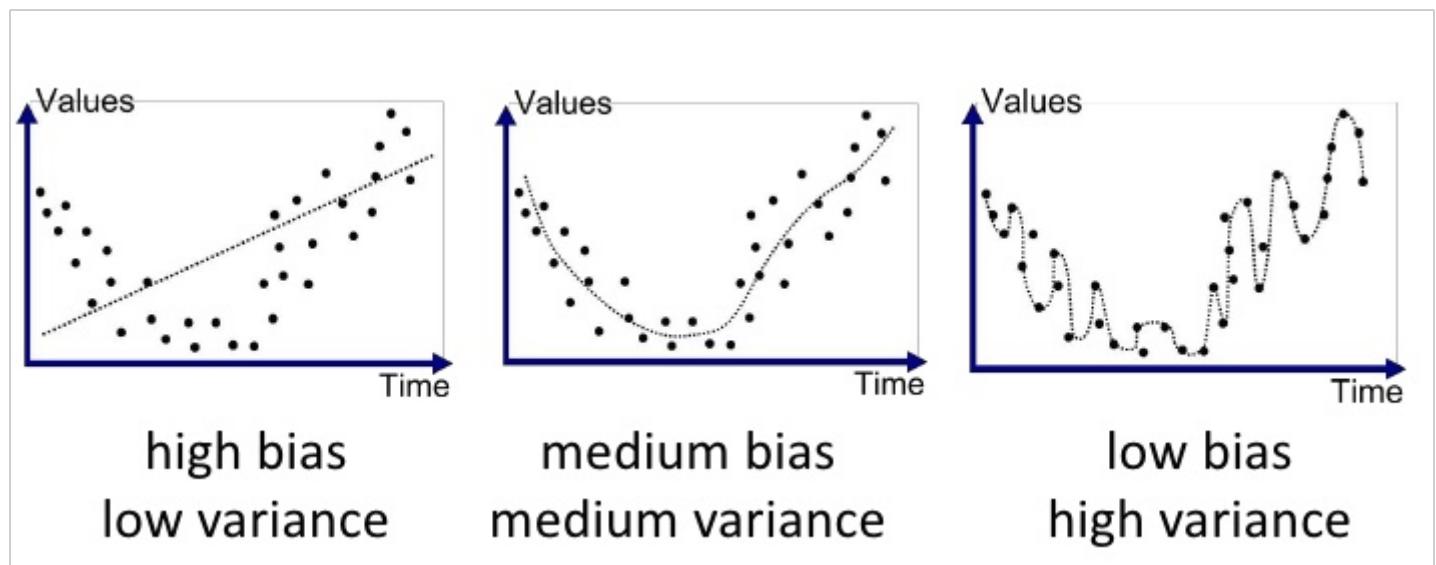
Evaluation of our method: the bias variance problem

We have to simultaneously minimize two different errors

The bias is error from erroneous assumptions in the learning algorithm. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting)

The variance is error from sensitivity to small fluctuations in the training set. High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

If we fit our data perfectly on the training set, it might not be general. Our aim is to have a general model



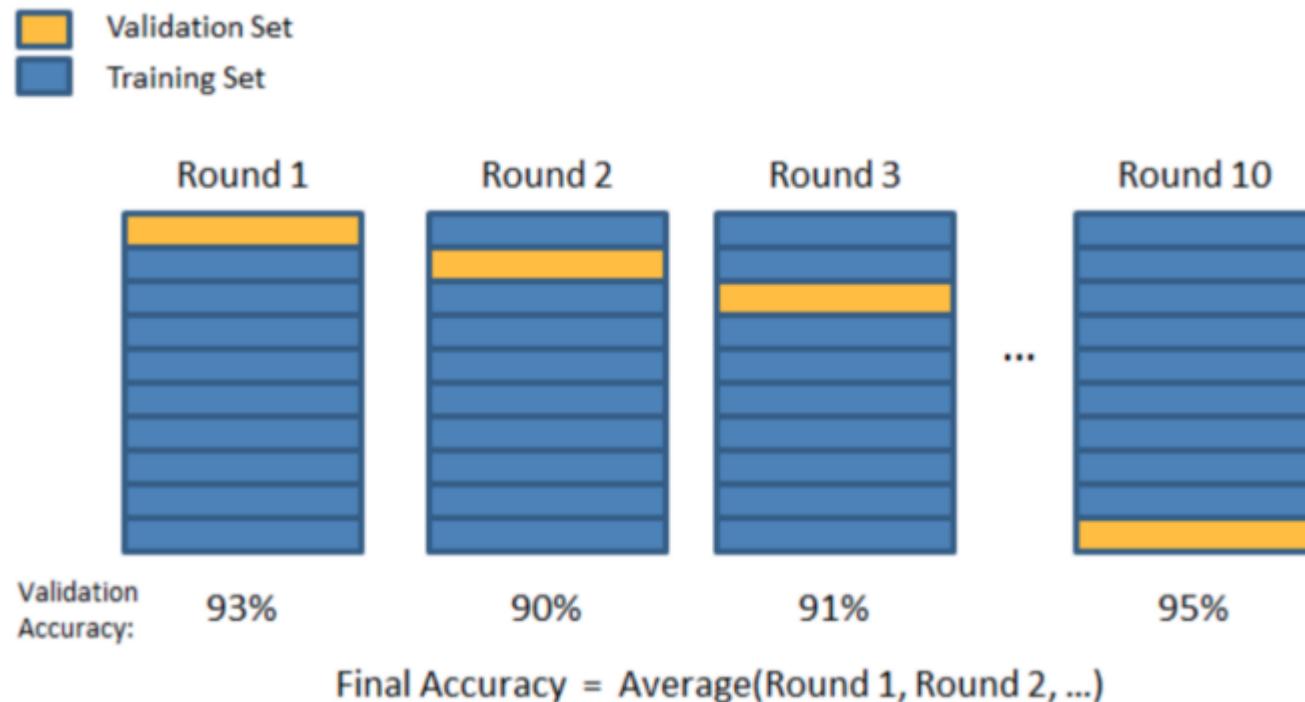
Evaluation of our method: the bias variance problem

We have to simultaneously minimize two different errors

The bias is error from erroneous assumptions in the learning algorithm. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting)

The variance is error from sensitivity to small fluctuations in the training set. High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

If we fit our data perfectly on the training set, it might not be general. Our aim is to have a general model



Evaluation of our method: the bias variance problem

We have to simultaneously minimize two different errors

The bias is error from erroneous assumptions in the learning algorithm. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting)

The variance is error from sensitivity to small fluctuations in the training set. High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

If we fit our data perfectly on the training set, it might not be general. Our aim is to have a general model

Let's reserve some of our data to see how general our model is:

```
1 #create train and test data
2 from sklearn.cross_validation import train_test_split
3 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.20, random_state = 5)
4 print(X_train.shape)
5 print(X_test.shape)
6 print(Y_train.shape)
7 print(Y_test.shape)
```

Once we have a good model we should re-run the training and evaluation on differently divided dataset and the accuracy should be consistent. This is called cross-validation

Scaling features

The features in our data can have orders of different magnitudes

Some algorithms can handle this fact, some can not

Even those which can learn these differences work faster and better with scaled data

Let's change all features to ensure that all of them has zero mean and variance of 1

```
1 import numpy as np
2 #scaling features
3 X_train=X_train-np.mean(X_train,0)
4 X_train=X_train/np.var(X_train,0)
5 X_test=X_test-np.mean(X_train,0)
6 X_test=X_test/np.var(X_train,0)
7 #from sklearn import preprocessing
8 #scaler = preprocessing.StandardScaler( )
9 #X_train = scaler.fit_transform(X_train)
```

Detecting outliers

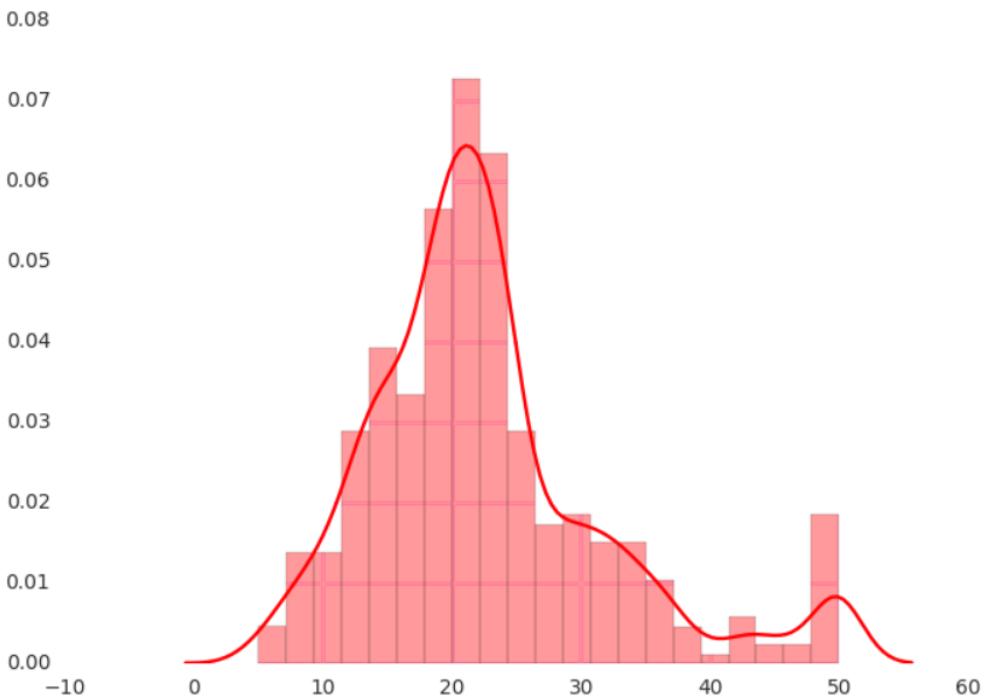
It really helps to have a look at the data:

```
1 print(X_test.shape)
2 fig,ax= plt.subplots()
3 sns.distplot(Y_train, color = 'red')
4 display(fig)
```

One could see that prices are limited at 50K\$

These elements will not fit into our model

They are outliers and it is not good to test the model on them in real life



The values are also thresholded in the test data in this case, so it is good if the model could learn this fact

Linear regression

We assume that there is a correlation in our data (x and y) which can be described as:
 $Y=wX+b$

We want to minimize the error between:

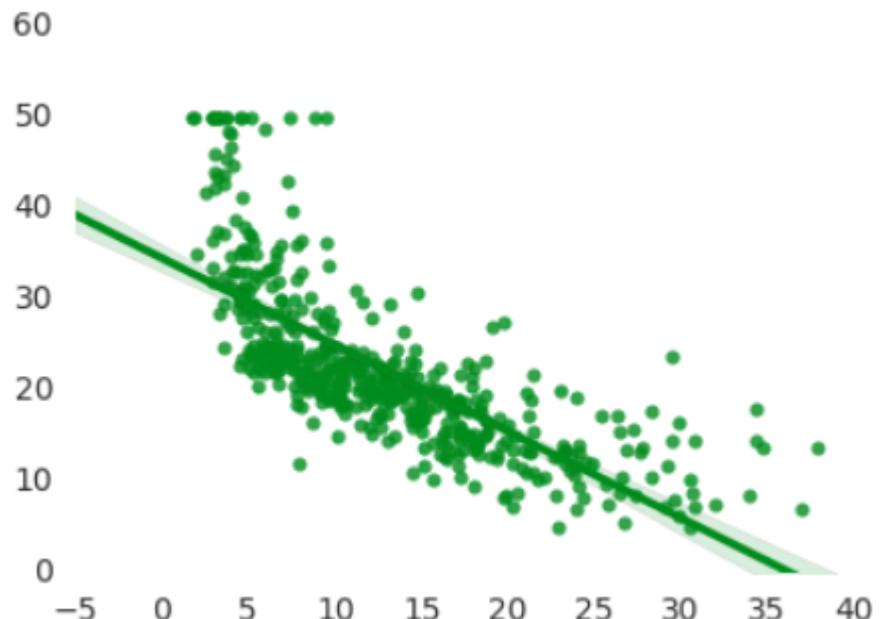
$$\sum_i |(wX_i + b) - Y_i|$$

According to w and b

There is an analytical formula to find the best w and b

We have to find the local extremum of the error function

$$\frac{\partial}{\partial w} \sum_i |(wX_i + b) - Y_i| = 0$$



$$\frac{\partial}{\partial b} \sum_i |(wX_i + b) - Y_i| = 0$$

Linear regression

We assume that there is a correlation in our data (x and y) which can be described as:
 $Y = wX + b$

Linear regression in SK learn

```
1 #use linear regression- we will fit the following model:
2 # y = wx + b
3 from sklearn.linear_model import LinearRegression
4
5 lm = LinearRegression()
6 lm.fit(X_train, Y_train)
7
8 Y_pred = lm.predict(X_test)
9
10 mse = sklearn.metrics.mean_squared_error(Y_test, Y_pred)
11 fig,ax= plt.subplots()
12
13 plt.scatter(Y_test, Y_pred)
14 plt.xlabel("Prices: $Y_i$")
15 plt.ylabel("Predicted prices: $\hat{Y}_i$")
16 plt.title("Prices vs Predicted prices, error: "+str(mse))
17 display(fig)
```

Linear regression

We assume that there is a correlation in our data (x and y) which can be described as:
 $Y=wX+b$

Linear regression in SK learn

```
1 #crossvalidation with simple linear regression
2 from sklearn.linear_model import LinearRegression
3 from sklearn. cross_validation import cross_val_score
4 linear = LinearRegression()
5 linear_scores = cross_val_score(linear, boston.data, boston.target, cv=5, scoring='neg_mean_squared_error')
6 print(linear_scores)
7 print(linear_scores.mean())
```

Linear regression - Lasso

We assume that there is a correlation in our data (x and y) which can be described as:
 $Y=wX+b$

One can easily overfit the data:

There is noise on our input data (some variables might be more reliable than the other)

How could we select the important variables.

Lasso (Least Absolute Shrinkage and Selection Operator) penalizes large values in the regressor (the same output accuracy where the values in w are more uniform is better)

The parameter of the algorithm is the penalty constant

```
1 #cross validation with lasso regression
2 from sklearn.linear_model import Lasso
3 lasso = Lasso(0.3) #score 0.1 to 0.5
4 lasso_scores = cross_val_score(lasso, boston.data, boston.target, cv=5, scoring='neg_mean_squared_error')
5 print(lasso_scores.mean())
6 print(lasso_scores)
```

Linear regression – Ridge, Elastic Net

We assume that there is a correlation in our data (x and y) which can be described as:
 $Y = wX + b$

One can easily overfit the data:

There is noise on our input data (some variables might be more reliable than the other)

Ridge and ElasticNet uses more complex penalization on the parameters depending on the data (Tikhonov regularization)

```
1 #cross validation with Ridge regression
2 from sklearn.linear_model import Ridge
3 ridge = Ridge(1) #score 1 to 5
4 ridge_scores = cross_val_score(ridge, boston.data, boston.target, cv = 5, scoring='neg_mean_squared_error')
5 print(ridge_scores.mean())
6 print(ridge_scores)
```

```
1 #cross validation using ElasticNet
2 from sklearn.linear_model import ElasticNet
3 elastic_net = ElasticNet(1) #score 1 to 5
4 elastic_net_scores = cross_val_score(elastic_net, boston.data, boston.target, cv = 5, scoring='neg_mean_absolute_error') #'neg_mean_absolute_error'
5 print(elastic_net_scores.mean())
6 print(elastic_net_scores)
```

Non-linear regression

How could we fit a more complex mode?

We could fit any model we found according to the previous formula:

$$Y = m(x)$$

Where we can derivate the model according to all the parameters and we will have an equation for every parameter

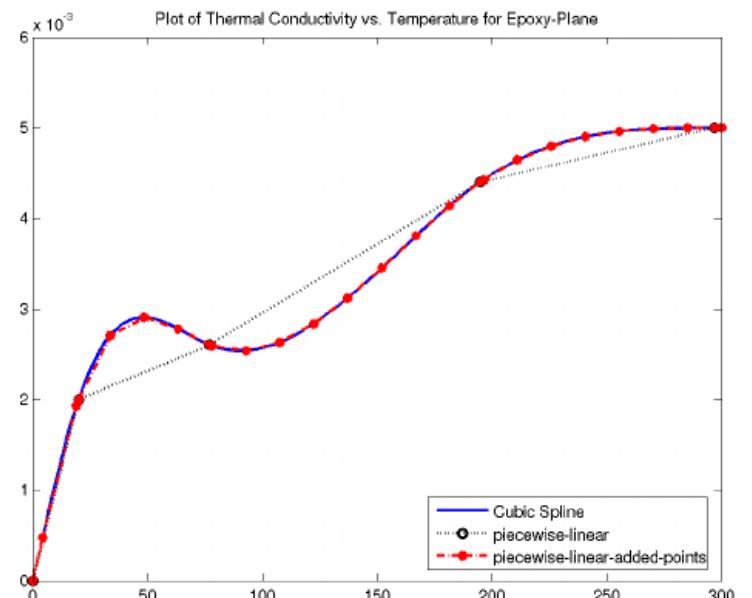
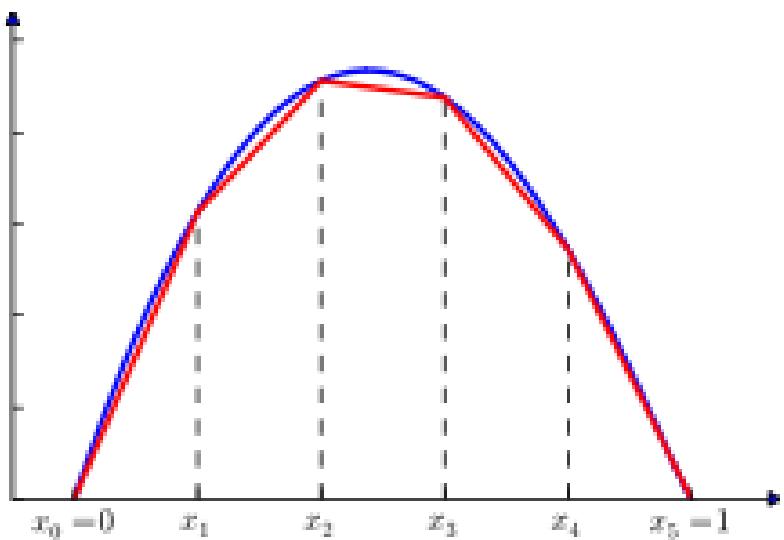
$$\frac{\partial}{\partial m_p} \sum_i |m(X_i) - Y_i| = 0$$

We could fit any model if it is known, but unfortunately in practice the problem is that we do not know the model

How could we create a general model, which could approximate any possible models

Non-linear regression

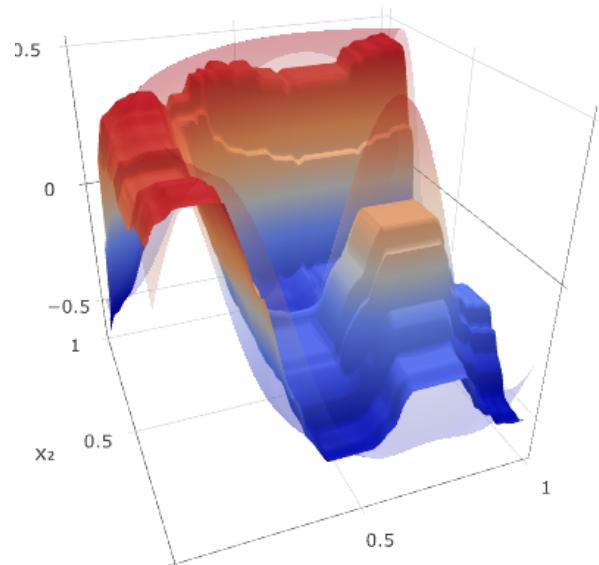
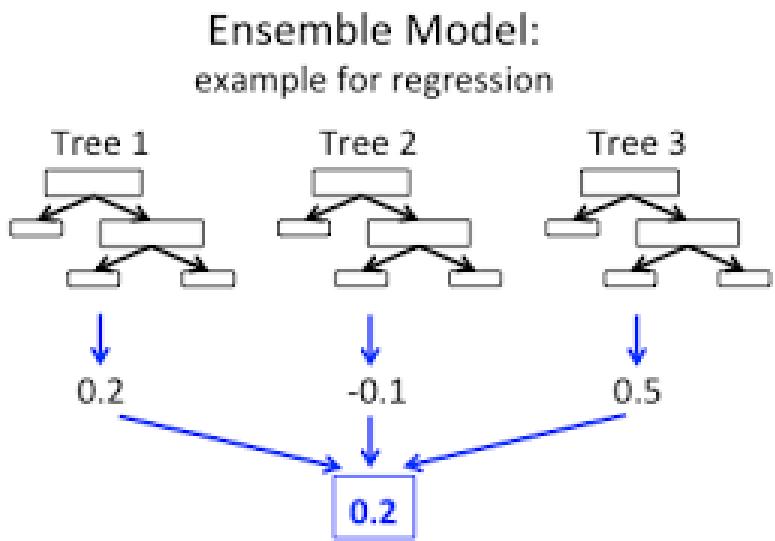
We can approximate any function with partially linear functions



We can divide the problem into sub-domains depending on the input values and use separated (or combined) linear regressors to approximate the original function.

Ensemble regression

We can approximate any function with partially linear functions



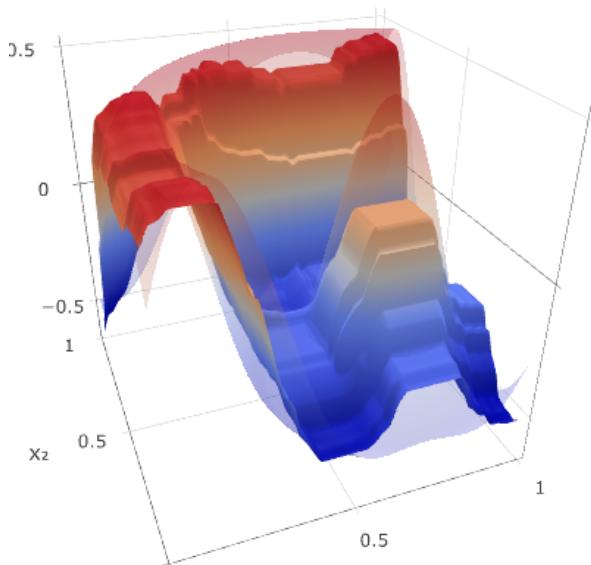
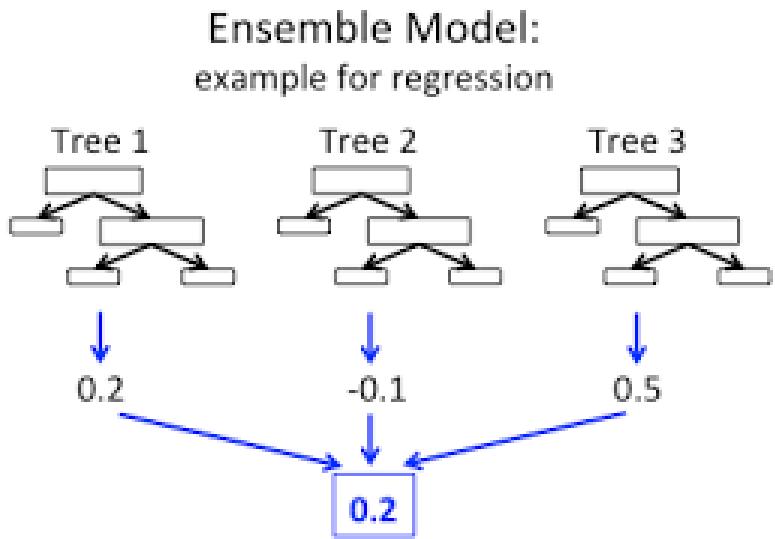
The approximation will be the weighted summation of linear models

The weights are determined by the values of the data

http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html

Ensemble regression

We can approximate any function with partially linear functions



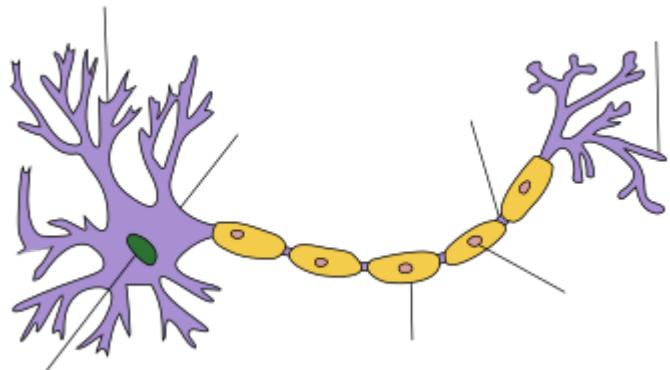
```
1 from sklearn import ensemble
2 ensemb = ensemble.GradientBoostingRegressor(
3     n_estimators=100, #number of independent estimators
4     learning_rate=0.1,
5     max_depth=4, #depth of the decision tree
6     min_samples_leaf=5, #minimum number of samples required to be a leaf node
7     max_features=0.8, #maximum numbers of features considered at each split
8     loss= 'lad', #least absolute deviation
9     random_state=700
10 )
11 ensemble_scores = cross_val_score(ensemb, boston.data, boston.target, cv = 5, scoring='neg_mean_absolute_error' ) #'neg_mean_absolute_error'
12 print(ensemble_scores.mean())
13 print(ensemble_scores)
```

Neural networks

Linear regression:

$$C_1(wx+b) + C_2(wx+b) + C_3(wx+b) + C_4(wx+b) + C_5(wx+b) + C_6(wx+b)$$

What does a neuron do?

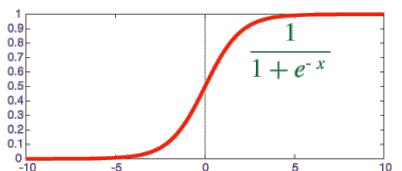
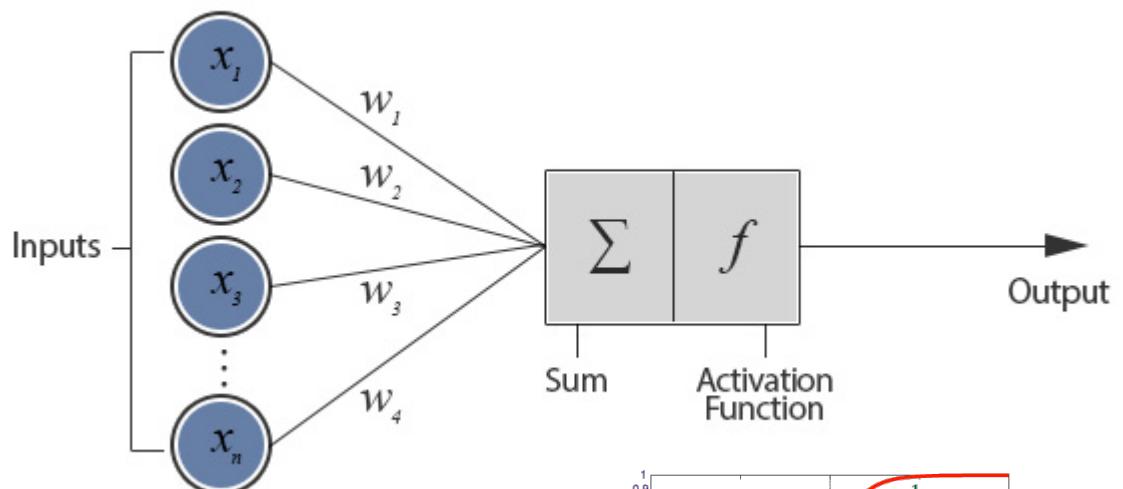
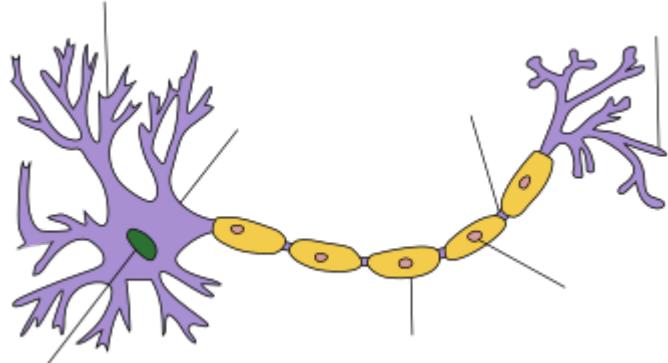


Neural networks

Linear regression:

$$C_1(wx+b) + C_2(wx+b) + C_3(wx+b) + C_4(wx+b) + C_5(wx+b) + C_6(wx+b)$$

What does a neuron do?



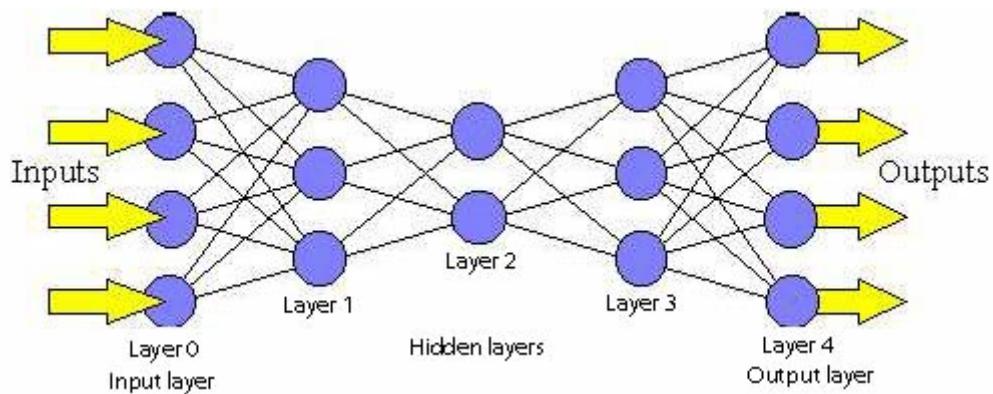
Neural networks

Linear regression:

$$C_1(wx+b) + C_2(wx+b) + C_3(wx+b) + C_4(wx+b) + C_5(wx+b) + C_6(wx+b)$$

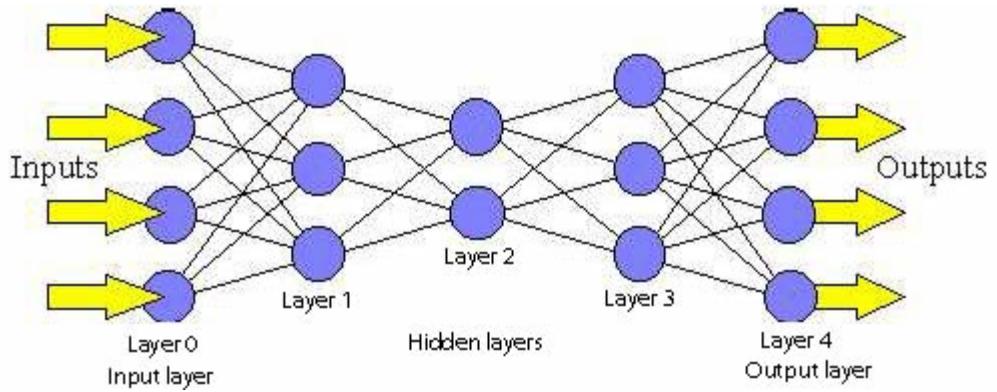
A feed forward neural network – layers of neurons:

$$f((w_3f((w_2f((w_1x)+b_1))+b_2))+b_3)$$



Neural networks

How can we define this as mathematical operations:

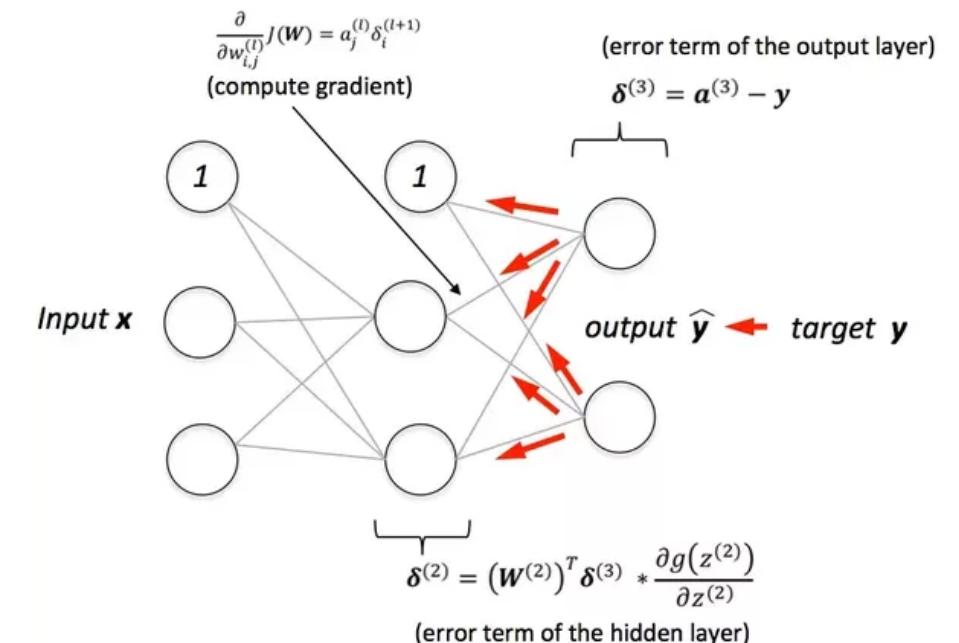
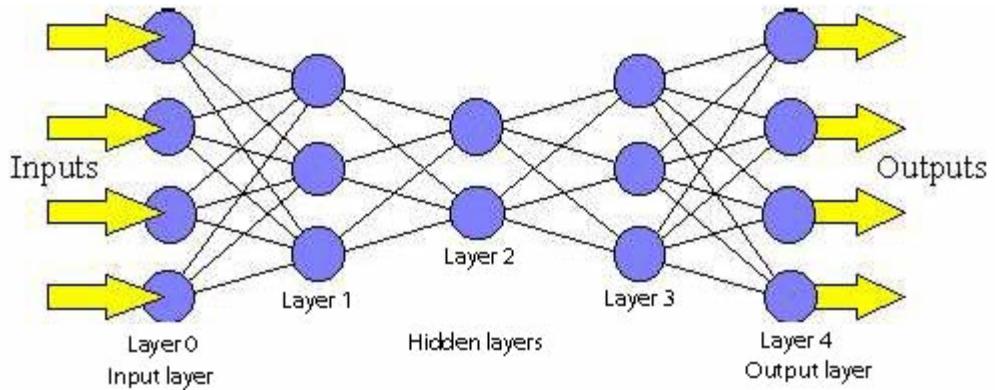


Neural networks

A feed-forward network is great

How could we teach such a network:

- Backpropagation algorithm

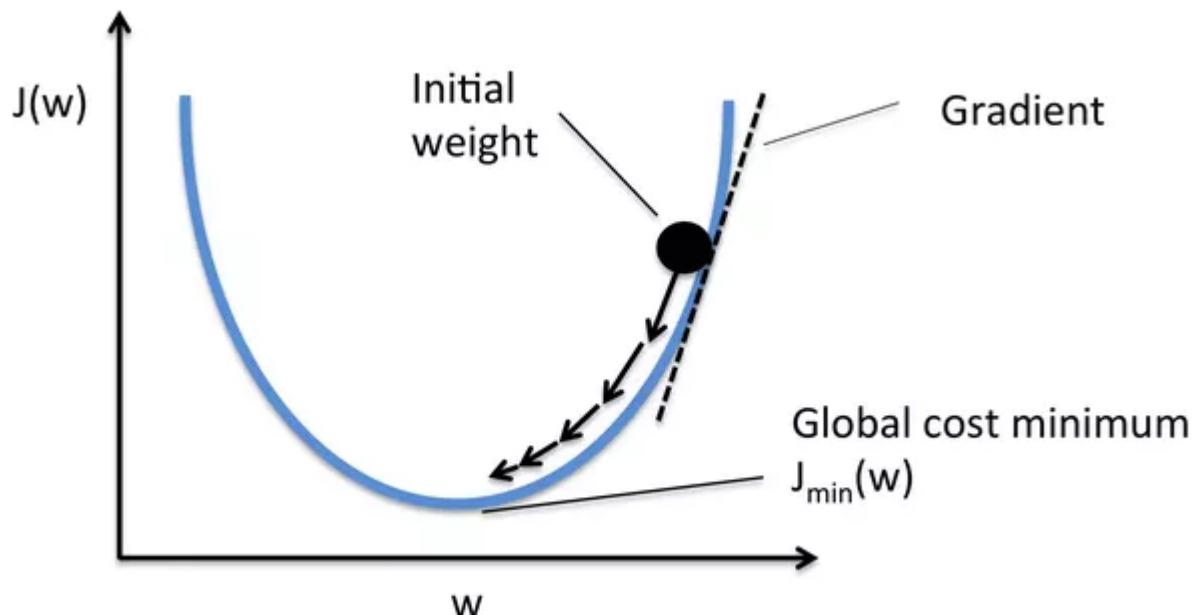
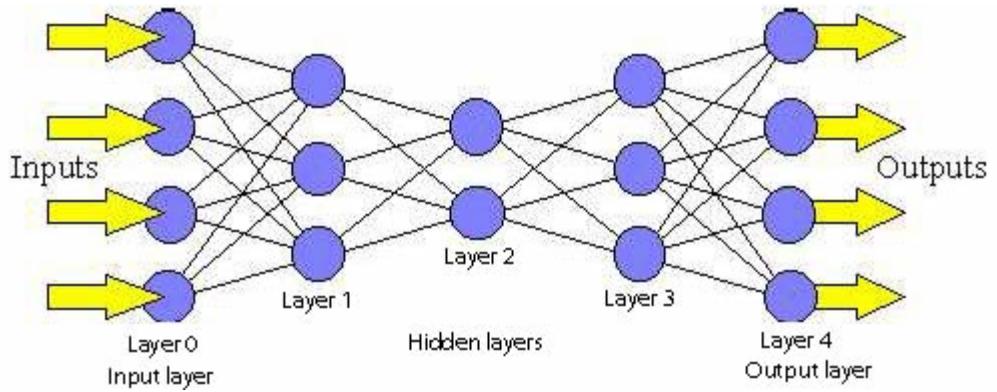


Neural networks

A feed-forward network is great

How could we teach such a network:

- Backpropagation algorithm
- Stochastic gradient descent



Neural networks

A feed-forward network is great

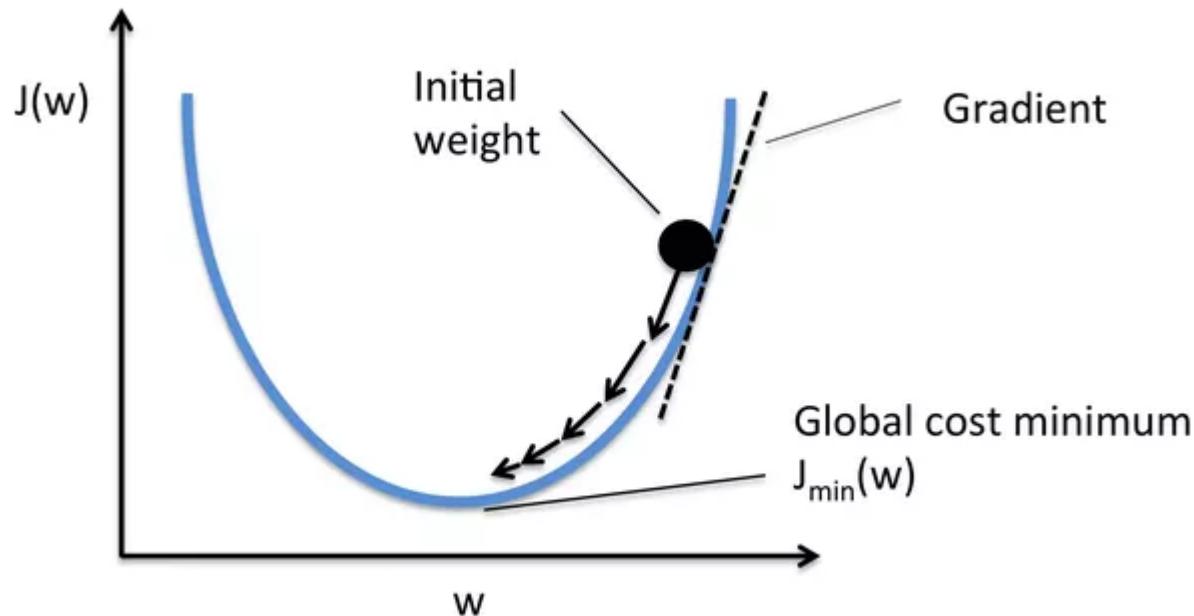
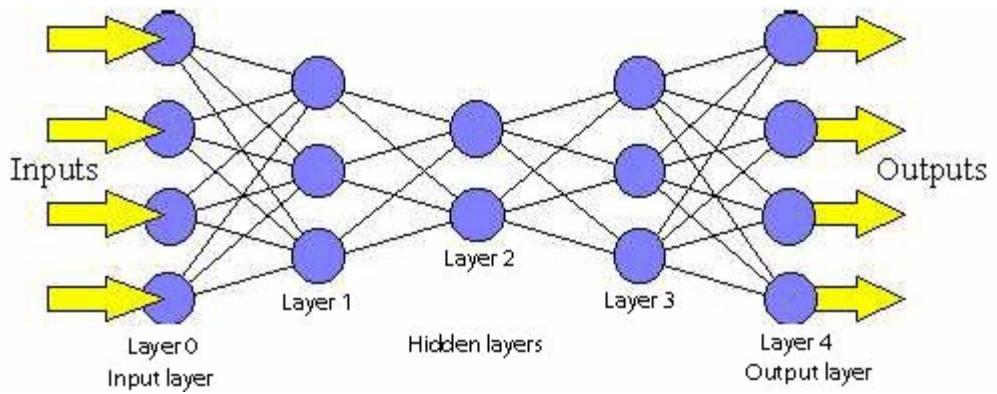
How could we teach such a network:

- Backpropagation algorithm
- Stochastic gradient descent

Tensorflow is here to help



TensorFlow

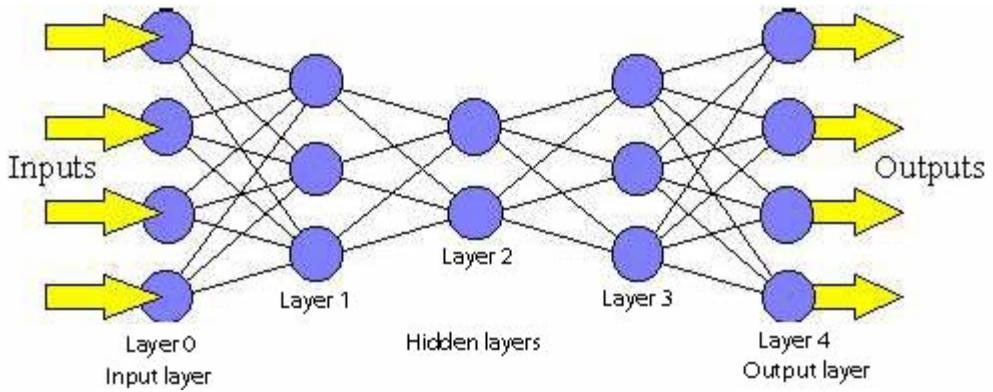


Neural networks

Regularization:

Using batches

Adding Dropout



Neural networks

```
12 tf.reset_default_graph()
13 # Parameters
14 learning_rate = 0.1
15 num_steps = 30000 #number of steps to train the network - one can stop training earlier
16 BatchLength = 52 #number of data fed paralell to the network
17
18 # Network Parameters
19 NumNeurons=[13,64,64,64,64,64,1] #the number of neruons in each layer
20
21 # tf Graph input -placeholders to feed data to the network
22 InputData = tf.placeholder(tf.float32, [None, NumNeurons[0]])
23 DesiredOutput = tf.placeholder(tf.float32, [None,1])
24 KeepProb = tf.placeholder(tf.float32)
25
26 CurrentInput=InputData
27 #this loop wil generate our network layer by layer
28 for i in range(len(NumNeurons)-1):
29     with tf.variable_scope('layer'+str(i)):
30         W=tf.Variable(tf.random_normal([NumNeurons[i], NumNeurons[i+1]], 0, 0.1))
31         B=tf.Variable(tf.random_normal([NumNeurons[i+1]], 0, 0.1, ))
32         CurrentInput = tf.add(tf.matmul(CurrentInput, W), B)
33         if i!=(len(NumNeurons)-2):
34             #the nonlinearity of the network
35             CurrentInput=tf.sigmoid(CurrentInput)
36             #CurrentInput=tf.nn.relu(CurrentInput)
37             #random dropout to keep variance
38             CurrentInput= tf.nn.dropout(CurrentInput, KeepProb)
39 Output=CurrentInput
40
41 #definitaiton of the error of the network
42 loss = tf.reduce_mean( tf.abs( tf.subtract(DesiredOutput , Output)))
43 optimizer = tf.train.AdagradOptimizer(learning_rate).minimize(loss)
44 #optimizer = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
45 #optimizer = tf.train.AdamOptimizer(learning_rate).minimize(loss)
```