# Logistic Regression

Sigmoid Activation in Perceptron

Bayes

Perceptron

Linear regression
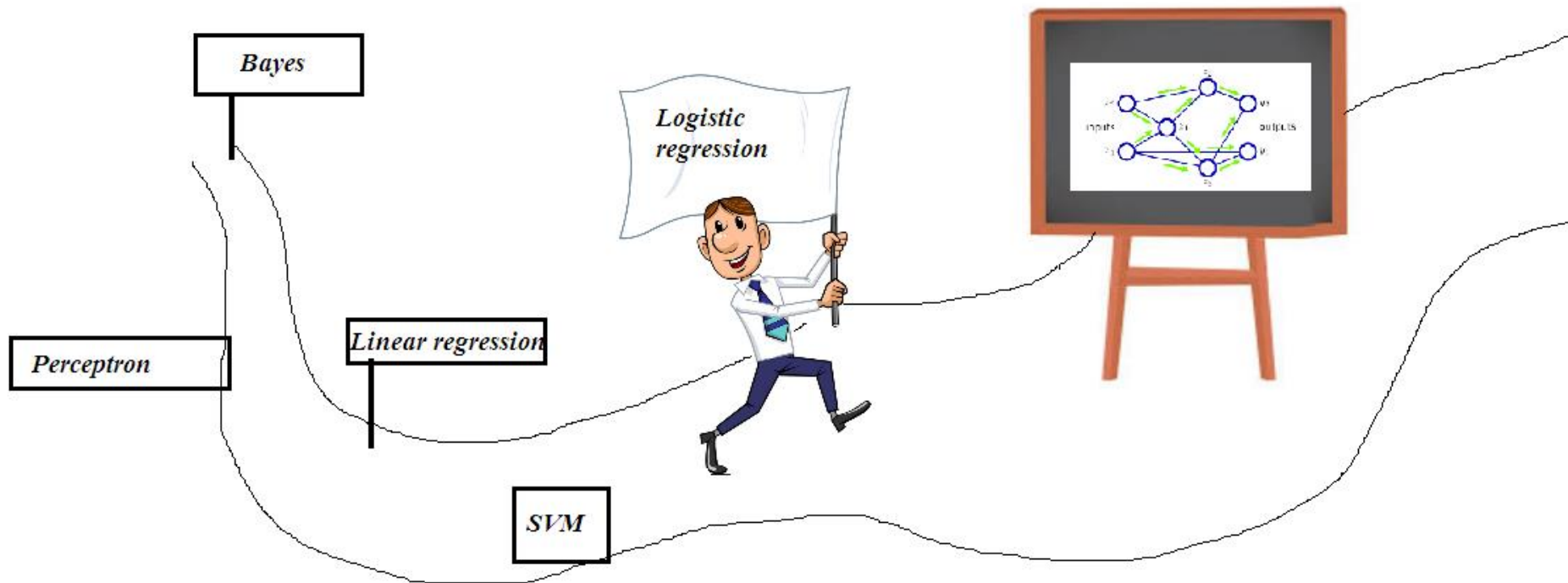
Logistic regression

SVM

inputs
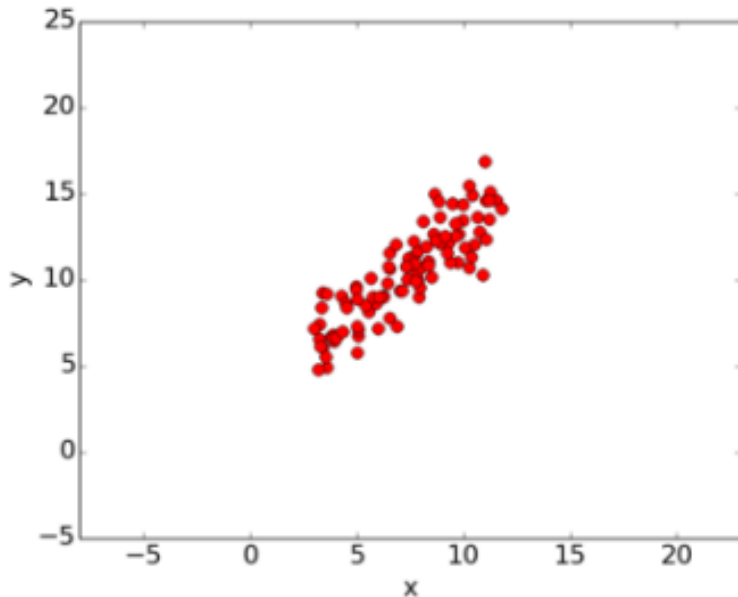
outputs

# Overview

- Linear Regression – Review
- Linear Regression – as classifier – Problem with outliers
- Logistic Regression – Sigmoid
  - Classifier ?
  - Criterion used in logistic regression
    - Why don't we use sum of squared error ?
  - Solution

# Linear regression



We want to fit a straight line (in 2D case). The sample we are given with is $\{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\}$.

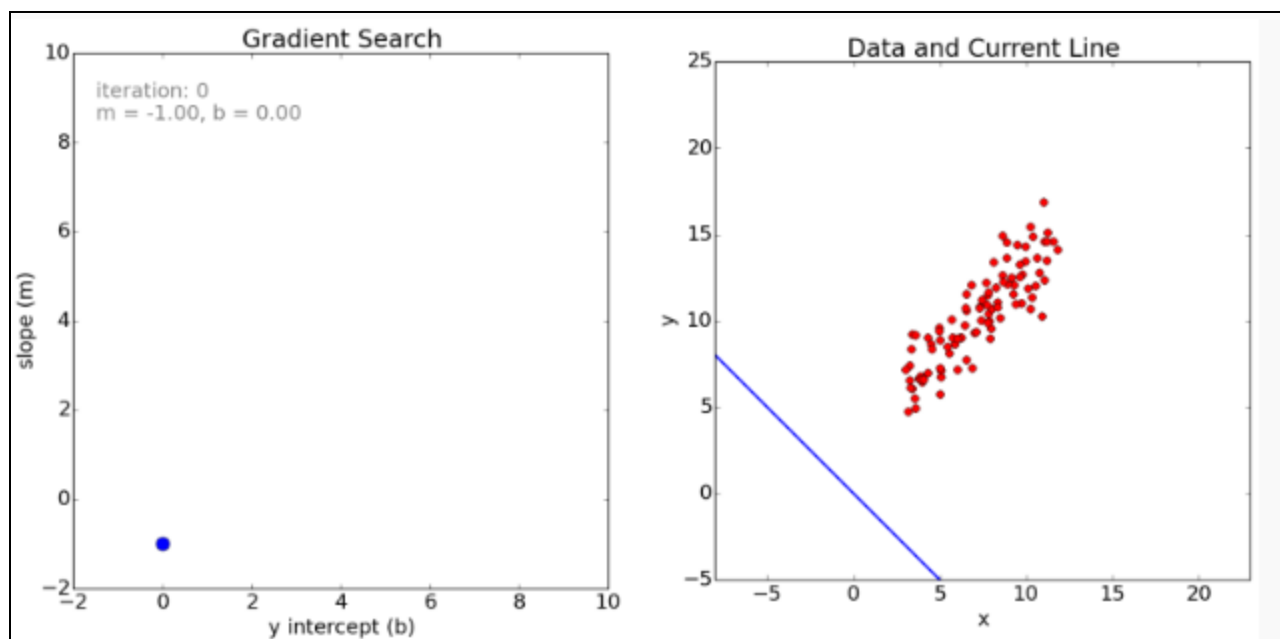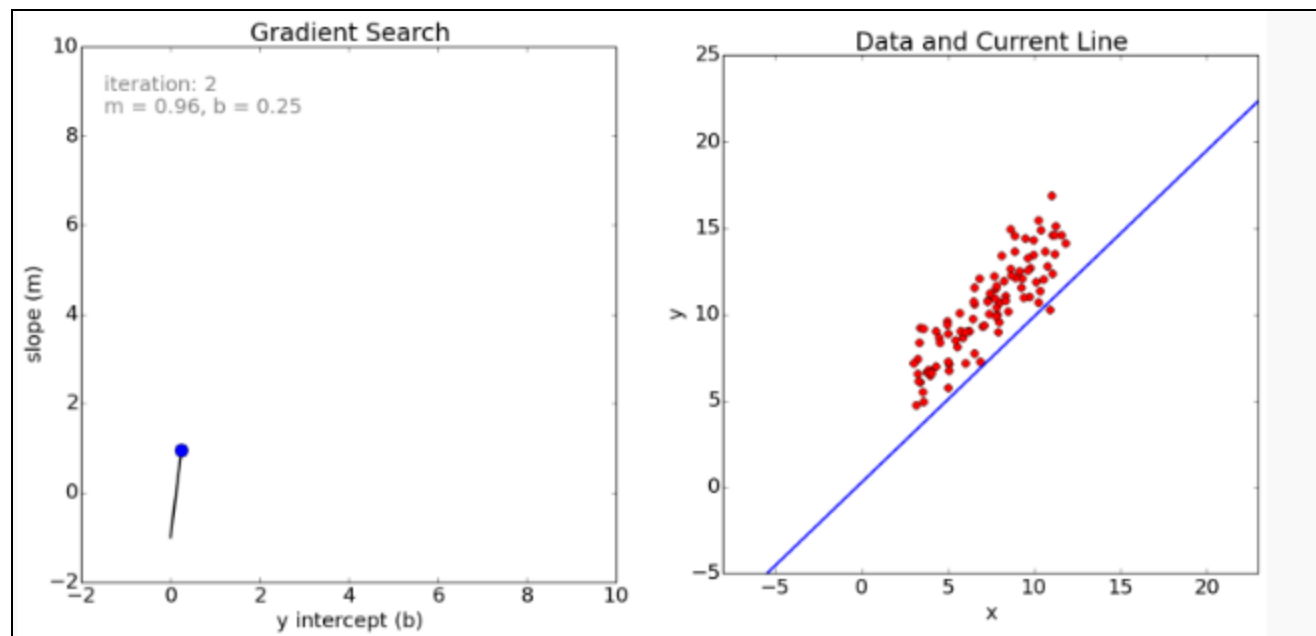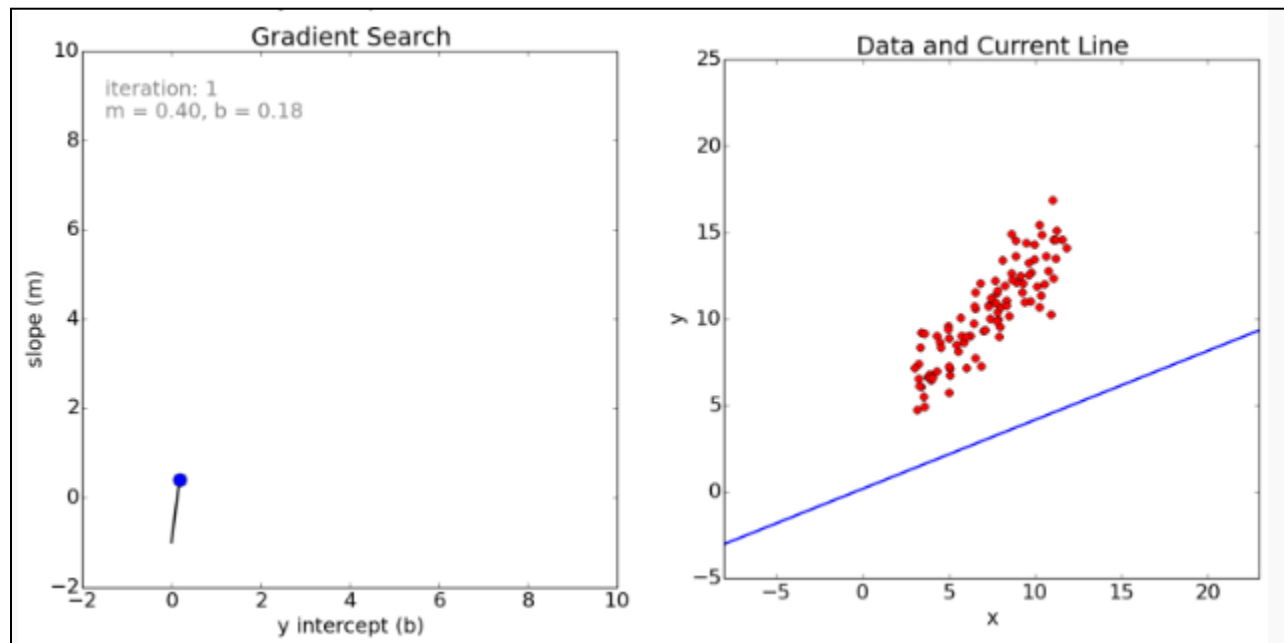We want to find (m, b).

In the space (m, b) what is the function we are going to define? Minimum value of that function should be a solution for us.
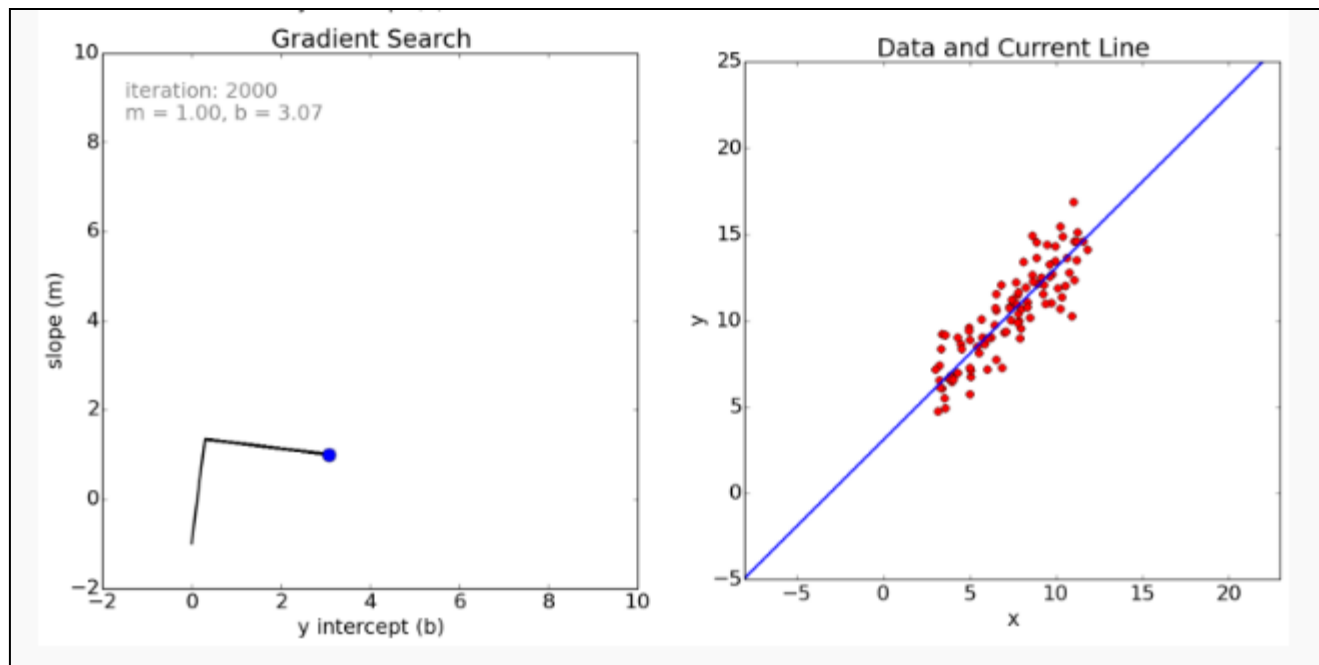
$$y = mx + b$$

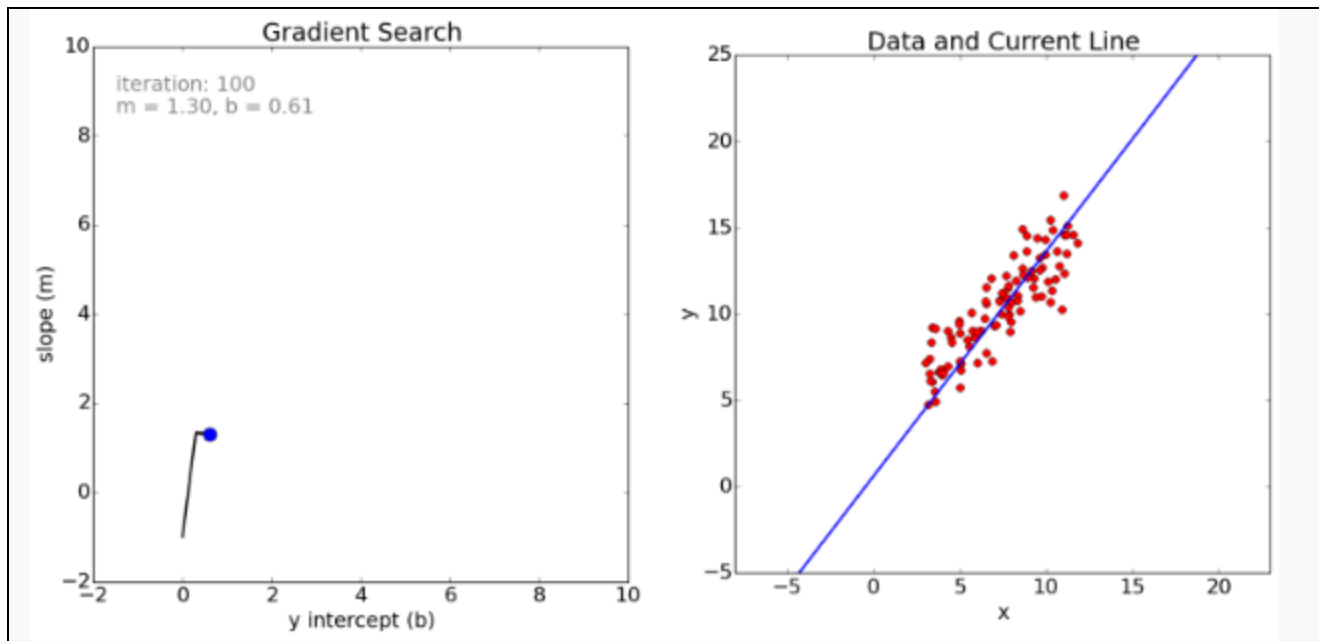$$\text{Error}_{(m,b)} = \frac{1}{N} \sum_{i=1}^{N} (y_i - (mx_i + b))^2$$

$$\frac{\partial}{\partial m} = \frac{2}{N} \sum_{i=1}^{N} -x_i(y_i - (mx_i + b))$$

$$\frac{\partial}{\partial b} = \frac{2}{N} \sum_{i=1}^{N} -(y_i - (mx_i + b))$$



5

Gradient Search — iteration: 1, m = 0.40, b = 0.18. Data and Current Line.

Gradient Search — iteration: 2, m = 0.96, b = 0.25. Data and Current Line.
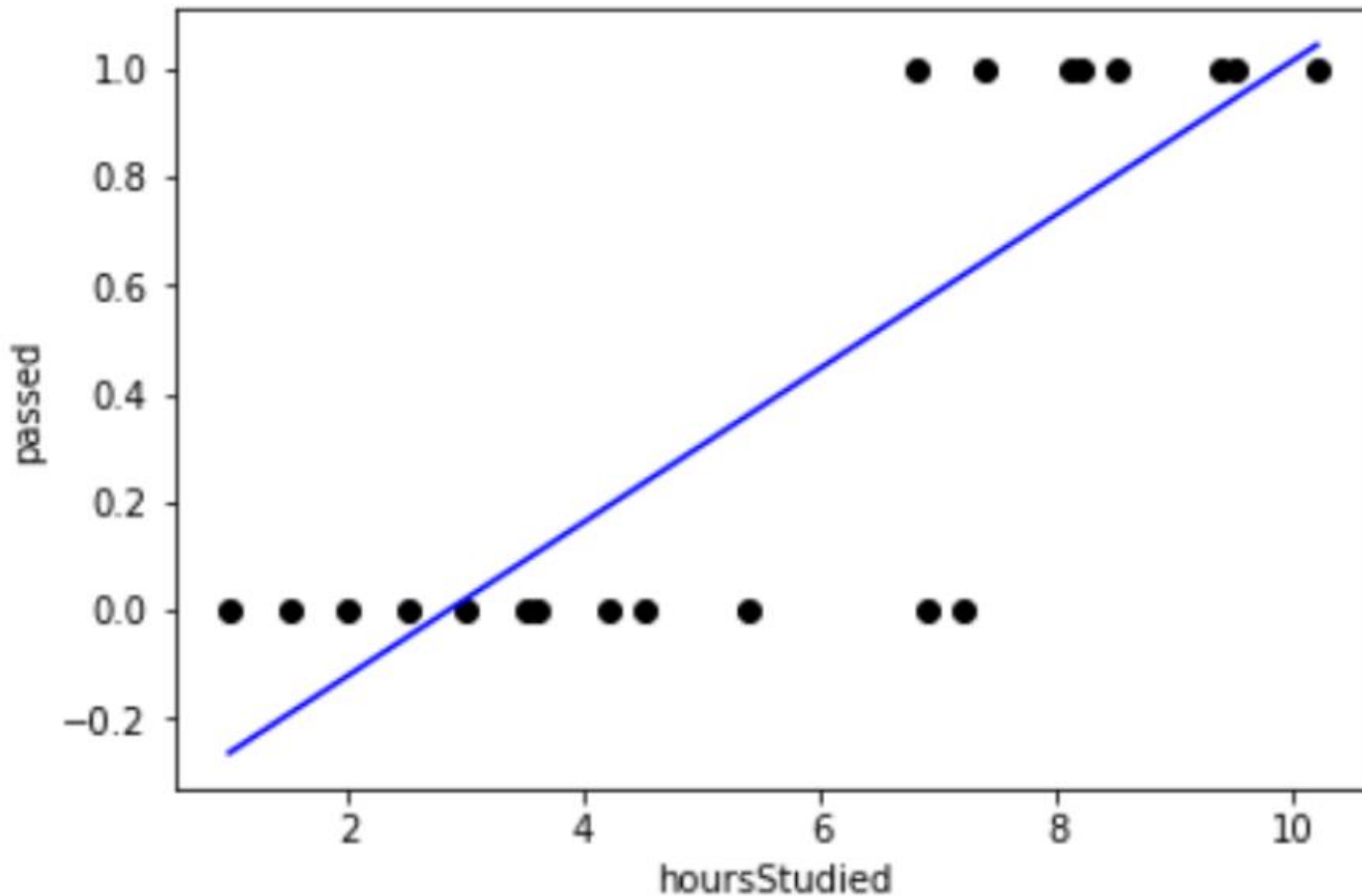
6

# Closed form solution

- Since the criterion that is minimized is quadratic, the linear regression problem must be having a closed form solution.

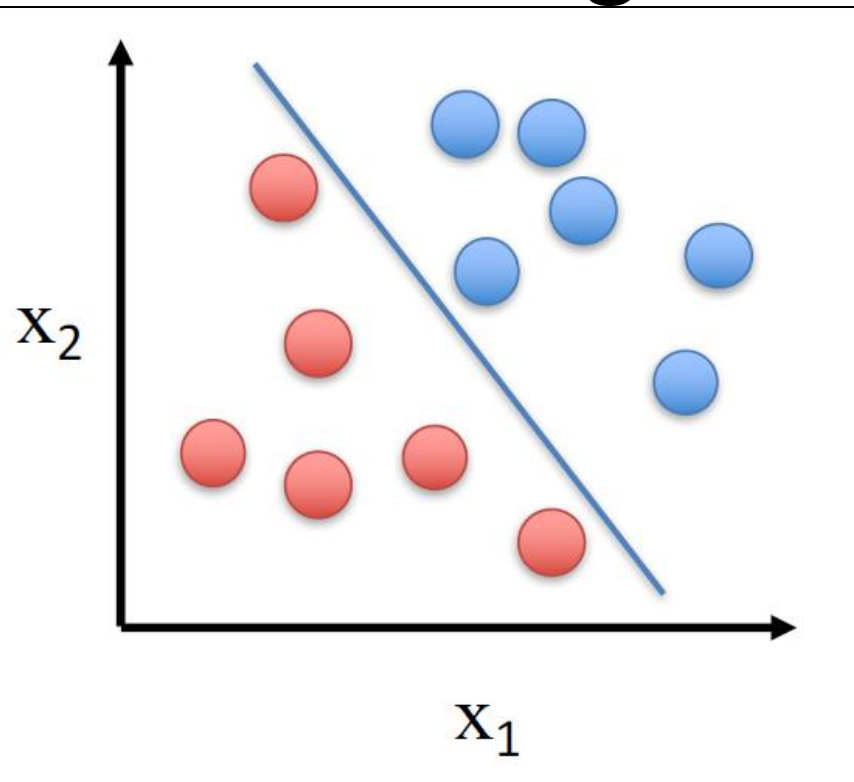- This is nothing but applying the Newton's descent method.

# Linear Regression for classification

Example with 1D data

# Linear Regression for classification



Example with 2D data

# Direct attempt, in learning the linear discriminant



$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \cdots \theta_d x_d$$

The parameter vector $\theta$ is learnt from the data such that the sum of squared error

$$E(\theta) = \frac{1}{2}\sum_{i=1}^{n}(y^{(i)} - h_\theta(x^{(i)}))^2$$

Where the training set

$E(\theta)\ is\ convex(also\ quadratic)$
$\Rightarrow no\ local\ minima\ problem\ (closed\ form\ solution)$

# Training Procedure

- $\nabla_\theta(E) = \sum_{i=1}^{n}\left(y^{(i)} - h_\theta\left(x^{(i)}\right)\right)\left(-x^{(i)}\right)$
- *Batch Method:*
- $\theta_{new} = \theta + \eta \sum_{i=1}^{n}\left(y^{(i)} - h_\theta\left(x^{(i)}\right)\right) x^{(i)}$
- Single Sample (Stochastic update)
- $\theta_{new} = \theta + \eta\left(y^{(i)} - h_\theta\left(x^{(i)}\right)\right)x^{(i)}$

----------------

- $h_\theta\left(x^{(i)}\right) = \theta^T x^{(i)}$
- $E(\theta) = \frac{1}{2}\sum_{i=1}^{n}\left(y^{(i)} - h_\theta\left(x^{(i)}\right)\right)^2$

# Linear Regression for classification



- But, outliers can be a big problem.

This is well within the blue class. This is a good training example.
**Why this is causing us problem??**

This is well within the blue class. This is a good training example.
**Why this is causing us problem??**

# Note

- Perceptron (Rosenblatt) does not have this problem (i.e., problem with seeming outliers).

# Note

- Perceptron (Rosenblatt) does not have this problem (i.e., problem with seeming outliers).

- But works only for linearly separable data.

# Note

- Perceptron (Rosenblatt) does not have this problem (i.e., problem with seeming outliers).
- But works only for linearly separable data.
- Actually, it is

- Step (or Sign) function is not differentiable.
- So we can't employ gradient descent to get the minimum error (in this formulation) solution.

# Logistic Regression

$x_0 = 1$

$x_1$   $\theta_1$   $\theta_0$

$x_2$   $\theta_2$

$\vdots$

$x_d$   $\theta_d$

$\Sigma$   $\sum_{i=0}^{d} \theta_i x_i$

$x_1$

$x_0 = 1$

$\theta_1$ $\theta_0$

$x_2$

$\theta_2$

$\vdots$

$\theta_d$

$x_d$

$\Sigma$

$\sum_{i=0}^{d} \theta_i x_i$

Sigmoid Activation Function

$\text{sig}(t) = \frac{1}{1+e^{-t}}$

$\text{sig}(t)$

1.0

0.8

0.6

0.4

0.2

$t$

−8    −6    −4    −2    2    4    6    8

# 1D : what it does



## Linear Regression



## Logistic Regression

# 1D : what it does



**Linear Regression**

Y=1

Y-Axis

Y=0

X-Axis

**Logistic Regression**



Logistic Regression

Category 1 observations
Category 2 observations

sigmoid(X)

X
(Measurement)

**See, the seeming outliers are indeed very good fits.**

# 2D data

# Notation used

- $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}, \; x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

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

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$$

Logistic / Sigmoid Function

$g(z)$

0.5

−6   −4   −2   0   2   4   6

# Notation used

- Given $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \left( \boldsymbol{x}^{(2)}, y^{(2)} \right), \ldots, \left( \boldsymbol{x}^{(n)}, y^{(n)} \right) \right\}$
  where $\boldsymbol{x}^{(i)} \in \mathbb{R}^d, \; y^{(i)} \in \{0, 1\}$

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}\right)$$

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

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}}}$$

Logistic / Sigmoid Function

$g(z)$

0.5

# Differentiation of the sigmoid

$$
\begin{aligned}
g'(z) &= \frac{d}{dz}\frac{1}{1+e^{-z}} \\[2mm]
&= \frac{1}{(1+e^{-z})^2}\left(e^{-z}\right) \\[2mm]
&= \frac{1}{(1+e^{-z})}\cdot\left(1-\frac{1}{(1+e^{-z})}\right) \\[2mm]
&= g(z)(1-g(z)).
\end{aligned}
$$

# Logistic Regression Objective Function

- Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$

  – Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\top} \boldsymbol{x}}}$$

  results in a non-convex optimization

# Classification Based on Probability

- Instead of just predicting the class, give the probability of the instance being that class
  - i.e., learn $p(y \mid \boldsymbol{x})$

- Comparison to perceptron:
  - Perceptron doesn't produce probability estimate

- Recall that:

$$0 \leq p(\text{event}) \leq 1$$

$$p(\text{event}) + p(\neg\text{event}) = 1$$

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

- $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should give $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

  – Want $0 \leq h_{\boldsymbol{\theta}}(\boldsymbol{x}) \leq 1$

Can't just use linear regression with a threshold

# Interpretation of Hypothesis Output

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$ = estimated  $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

Example:  Cancer diagnosis from tumor size

$$\boldsymbol{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that:  $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) + p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1$

Therefore,  $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1 - p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

# Logistic Regression

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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

$g(z)$

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>negative</u> values for negative instances

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>positive</u> values for positive instances

- Assume a threshold and…
  - Predict y = 1 if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5$
  - Predict y = 0 if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5$

y = 1

$\theta$

y = 0

# What we want …??

- We have

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

- Given



- Find $\theta$ so that the data agrees to the maximum extent.

- We can use maximum likelihood parameter estimation (MLE)

- Likelihood of data is given by: $l(\boldsymbol{\theta}) = \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$

- So, looking for the $\boldsymbol{\theta}$ that maximizes the likelihood

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{\theta}} \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

- Can take the log without changing the solution:

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} \log \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

Note that this can be written more compactly as

$$p(y \mid x; \theta) = (h_\theta(x))^y \, (1 - h_\theta(x))^{1-y}$$

# Likelihood

$$L(\theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$

$$= \prod_{i=1}^{n} \left(h_\theta(x^{(i)})\right)^{y^{(i)}} \left(1 - h_\theta(x^{(i)})\right)^{1-y^{(i)}}$$

- Log-likelihood

$$\ell(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^{n} y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))$$

# Find $\theta$ that maximizes $l(\theta)$

- $l(\theta)$ is taken to be the criterion, but, note this has to be maximized.

- We can employ gradient ascent method.

- We can show that, the negative of the log-likelihood, i.e., $-l(\theta)$ which should be minimized is convex (hence no local minima problem)

- For proof of this refer the link
  http://mathgotchas.blogspot.com/2011/10/why-is-error-function-minimized-in.html

- $\nabla_\theta l(\theta) = \big(y - h_\theta(x)\big)\, x$
  - Note, here $x$ is a vector.
  - This is for a single training example.

- This we obtained from,

$$
\begin{aligned}
\frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1-y)\frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\
&= \left( y \frac{1}{g(\theta^T x)} - (1-y)\frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\
&= \left( y(1 - g(\theta^T x)) - (1-y)g(\theta^T x) \right) x_j \\
&= \left( y - h_\theta(x) \right) x_j
\end{aligned}
$$

# Stochastic or single sample update rule

- $\theta_{new} \qquad = \theta + \eta \nabla_\theta l(\theta)$
$= \theta + \eta \big( y - h_\theta(x) \big) \, x$

## Batch Method

$$\theta_{new} = \theta + \eta \sum_{i=1}^{n} \big( y^{(i)} - h_\theta \big( x^{(i)} \big) \big) \, x^{(i)}$$

# Stochastic or single sample update rule

- $\theta_{new} \qquad = \theta + \eta \nabla_\theta l(\theta)$
$$= \theta + \eta \left( y - h_\theta(x) \right) x$$

## Batch Method

$$\theta_{new} = \theta + \eta \sum_{i=1}^{n} \left( y^{(i)} - h_\theta\left(x^{(i)}\right) \right) x^{(i)}$$

This looks IDENTICAL to linear regression!!!

- However, the form of the model is very different:

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^\top x}}$$

- The End

- Supplementary Material

- $h_\theta(x) = \sum_{i=1}^{d} \theta_i x_i$