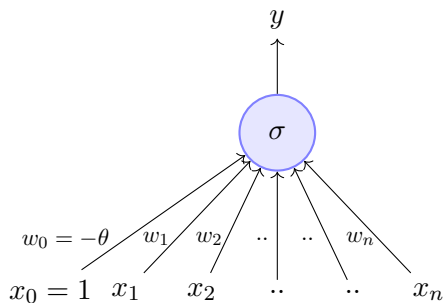
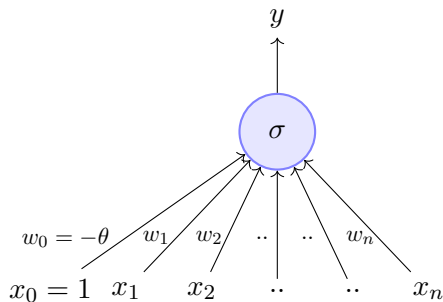


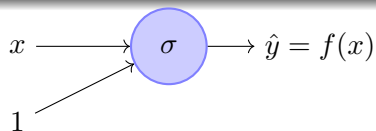
Module 3.3: Learning Parameters: (Infeasible) guess work



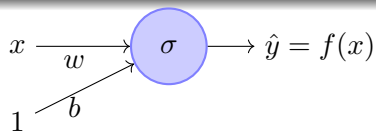
- With this setup in mind, we will now focus on this **model** and discuss an **algorithm** for learning the **parameters** of this model from some given **data** using an appropriate **objective function**



- With this setup in mind, we will now focus on this **model** and discuss an **algorithm** for learning the **parameters** of this model from some given **data** using an appropriate **objective function**
- σ stands for the sigmoid function (logistic function in this case)

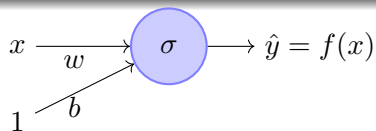


- With this setup in mind, we will now focus on this **model** and discuss an **algorithm** for learning the **parameters** of this model from some given **data** using an appropriate **objective function**
- σ stands for the sigmoid function (logistic function in this case)
- For ease of explanation, we will consider a very simplified version of the model having just 1 input



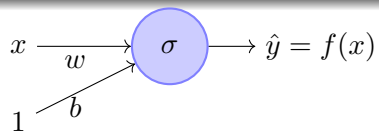
$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

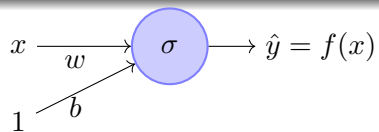
- With this setup in mind, we will now focus on this **model** and discuss an **algorithm** for learning the **parameters** of this model from some given **data** using an appropriate **objective function**
- σ stands for the sigmoid function (logistic function in this case)
- For ease of explanation, we will consider a very simplified version of the model having just 1 input
- Further to be consistent with the literature, from now on, we will refer to w_0 as b (bias)



$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

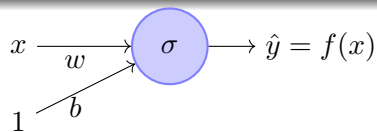
- With this setup in mind, we will now focus on this **model** and discuss an **algorithm** for learning the **parameters** of this model from some given **data** using an appropriate **objective function**
- σ stands for the sigmoid function (logistic function in this case)
- For ease of explanation, we will consider a very simplified version of the model having just 1 input
- Further to be consistent with the literature, from now on, we will refer to w_0 as b (bias)
- Lastly, instead of considering the problem of predicting like/dislike, we will assume that we want to predict *criticsRating(y)* given *imdbRating(x)* (for no particular reason)





Input for training

$\{x_i, y_i\}_{i=1}^N \rightarrow N$ pairs of (x, y)



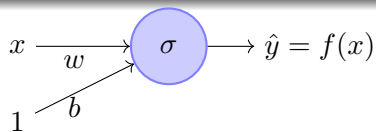
Input for training

$\{x_i, y_i\}_{i=1}^N \rightarrow N$ pairs of (x, y)

Training objective

Find w and b such that:

$$\underset{w, b}{\text{minimize}} \mathcal{L}(w, b) = \sum_{i=1}^N (y_i - f(x_i))^2$$



$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

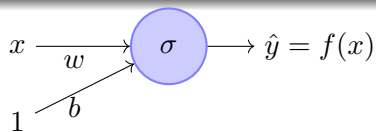
Input for training

$\{x_i, y_i\}_{i=1}^N \rightarrow N$ pairs of (x, y)

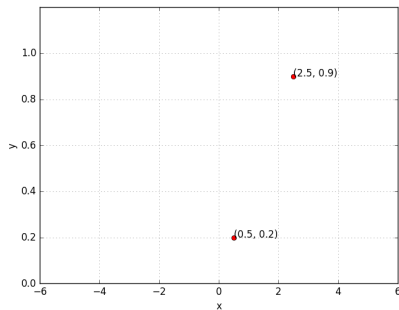
Training objective

Find w and b such that:

$$\underset{w, b}{\text{minimize}} \mathcal{L}(w, b) = \sum_{i=1}^N (y_i - f(x_i))^2$$

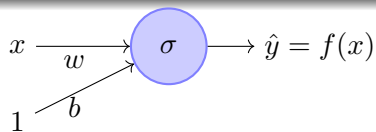


$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

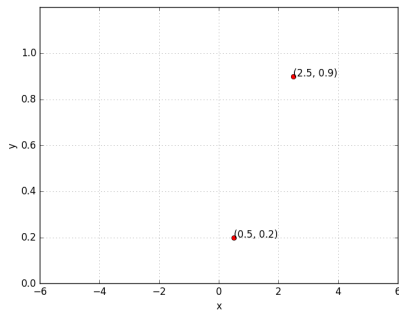


What does it mean to train the network?

- Suppose we train the network with $(x, y) = (0.5, 0.2)$ and $(2.5, 0.9)$

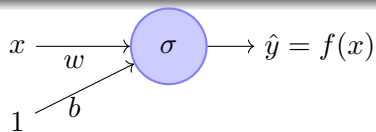


$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

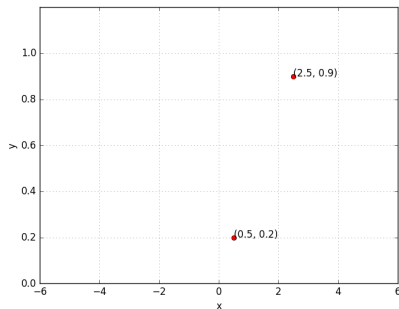


What does it mean to train the network?

- Suppose we train the network with $(x, y) = (0.5, 0.2)$ and $(2.5, 0.9)$
- At the end of training we expect to find w^* , b^* such that:

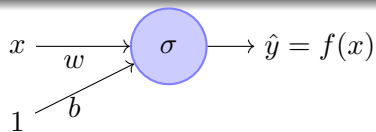


$$f(x) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

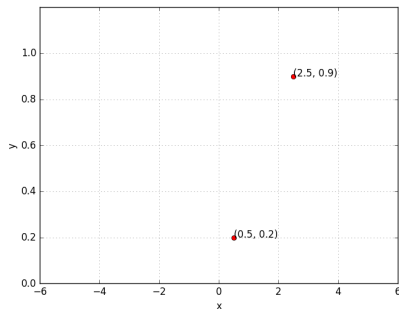


What does it mean to train the network?

- Suppose we train the network with $(x, y) = (0.5, 0.2)$ and $(2.5, 0.9)$
- At the end of training we expect to find w^* , b^* such that:
- $f(0.5) \rightarrow 0.2$ and $f(2.5) \rightarrow 0.9$



$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

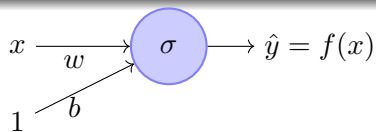


What does it mean to train the network?

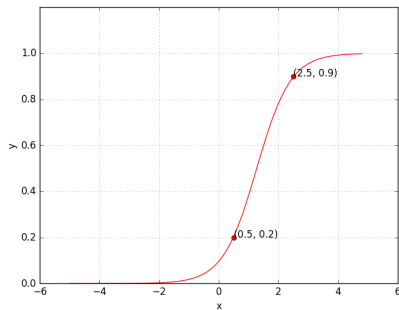
- Suppose we train the network with $(x, y) = (0.5, 0.2)$ and $(2.5, 0.9)$
- At the end of training we expect to find w^* , b^* such that:
- $f(0.5) \rightarrow 0.2$ and $f(2.5) \rightarrow 0.9$

In other words...

- We hope to find a sigmoid function such that $(0.5, 0.2)$ and $(2.5, 0.9)$ lie on this sigmoid



$$f(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$



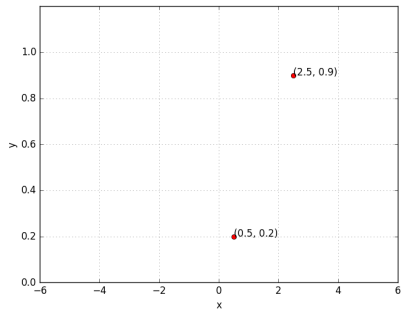
What does it mean to train the network?

- Suppose we train the network with $(x, y) = (0.5, 0.2)$ and $(2.5, 0.9)$
- At the end of training we expect to find w^* , b^* such that:
- $f(0.5) \rightarrow 0.2$ and $f(2.5) \rightarrow 0.9$

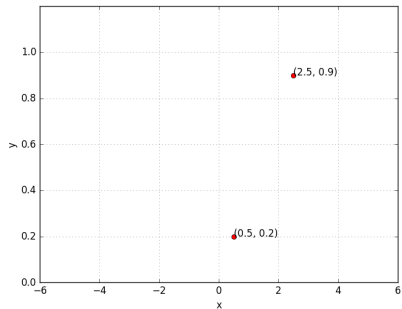
In other words...

- We hope to find a sigmoid function such that $(0.5, 0.2)$ and $(2.5, 0.9)$ lie on this sigmoid

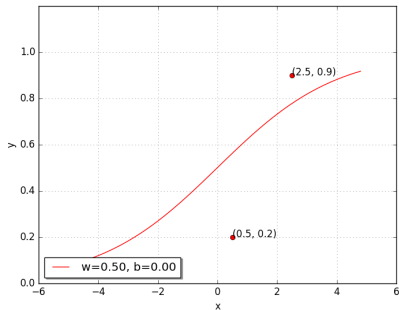
Let us see this in more detail....



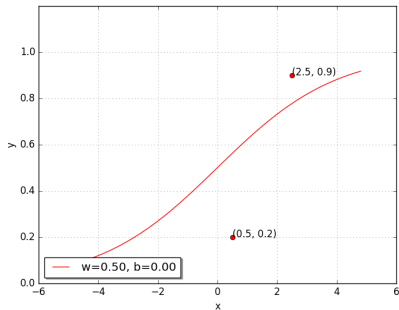
- Can we try to find such a w^*, b^* manually

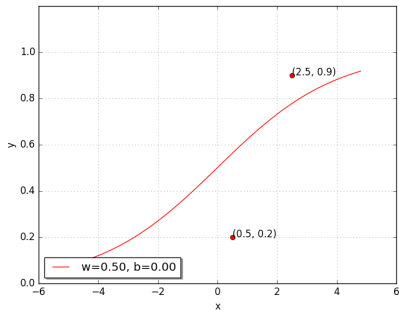


- Can we try to find such a w^*, b^* manually
- Let us try a random guess.. (say, $w = 0.5, b = 0$)

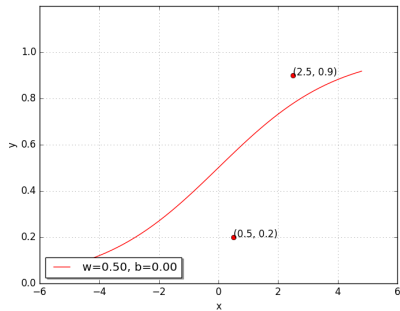


- Can we try to find such a w^*, b^* manually
- Let us try a random guess.. (say, $w = 0.5, b = 0$)
- Clearly not good, but how bad is it ?

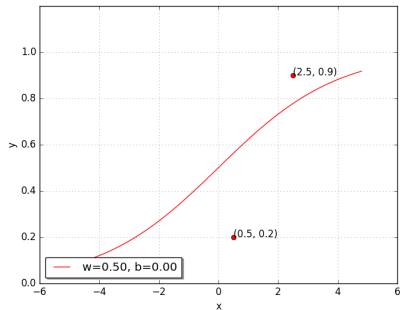




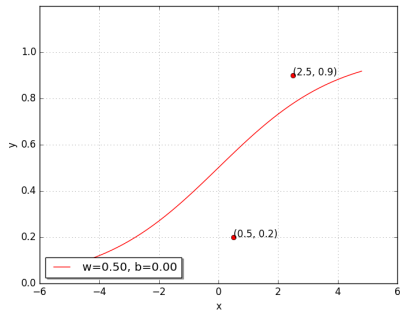
- Can we try to find such a w^*, b^* manually
- Let us try a random guess.. (say, $w = 0.5, b = 0$)
- Clearly not good, but how bad is it ?
- Let us revisit $\mathcal{L}(w, b)$ to see how bad it is ...



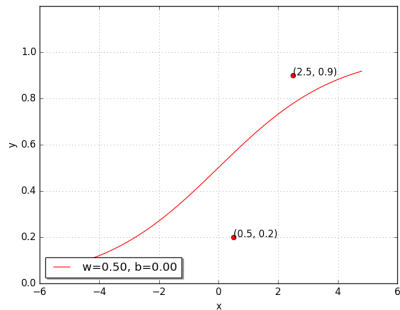
$$\mathcal{L}(w, b) = \frac{1}{2} * \sum_{i=1}^N (y_i - f(x_i))^2$$



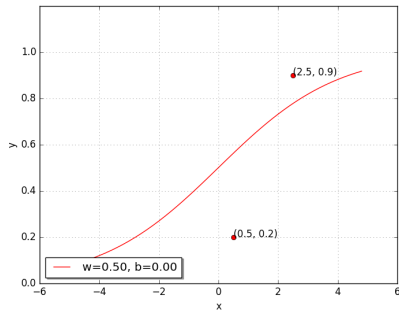
$$\begin{aligned}\mathcal{L}(w, b) &= \frac{1}{2} * \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \frac{1}{2} * (y_1 - f(x_1))^2 + (y_2 - f(x_2))^2\end{aligned}$$



$$\begin{aligned}\mathcal{L}(w, b) &= \frac{1}{2} * \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \frac{1}{2} * (y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 \\ &= \frac{1}{2} * (0.9 - f(2.5))^2 + (0.2 - f(0.5))^2\end{aligned}$$



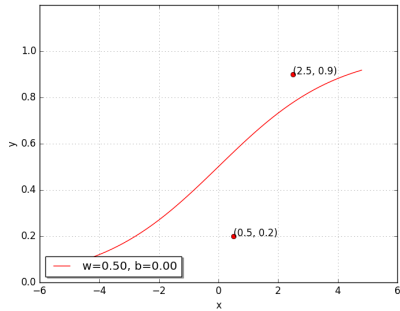
$$\begin{aligned}\mathcal{L}(w, b) &= \frac{1}{2} * \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \frac{1}{2} * (y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 \\ &= \frac{1}{2} * (0.9 - f(2.5))^2 + (0.2 - f(0.5))^2 \\ &= 0.073\end{aligned}$$



$$\begin{aligned}\mathcal{L}(w, b) &= \frac{1}{2} * \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \frac{1}{2} * (y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 \\ &= \frac{1}{2} * (0.9 - f(2.5))^2 + (0.2 - f(0.5))^2 \\ &= 0.073\end{aligned}$$

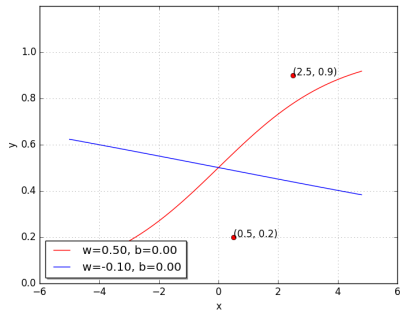
We want $\mathcal{L}(w, b)$ to be as close to 0 as possible

Let us try some other values of w , b



w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730

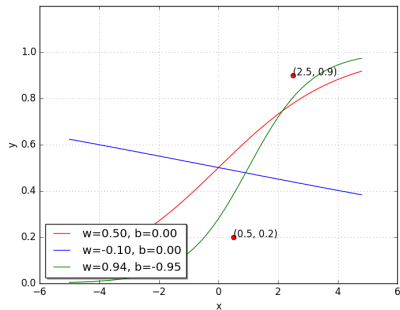
Let us try some other values of w , b



w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481

Oops!! this made things even worse...

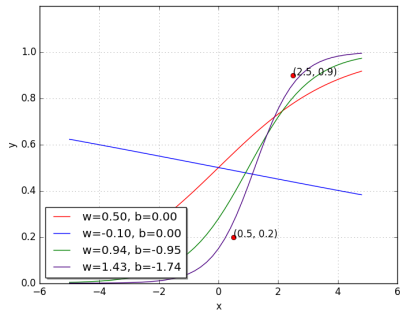
Let us try some other values of w , b



w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214

Perhaps it would help to push w and b in the other direction...

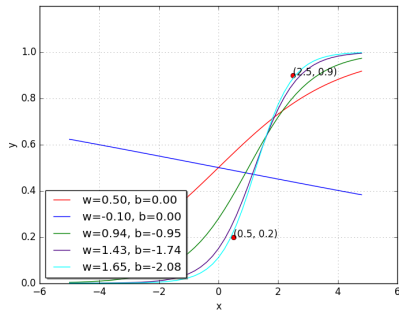
Let us try some other values of w , b



w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028

Let us keep going in this direction, *i.e.*, increase w and decrease b

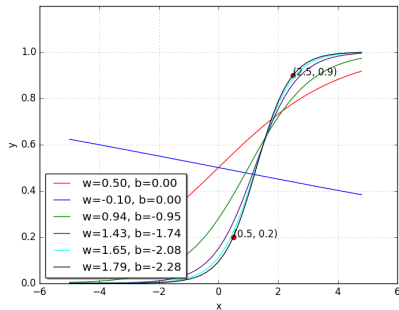
Let us try some other values of w , b



w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003

Let us keep going in this direction, *i.e.*, increase w and decrease b

Let us try some other values of w , b



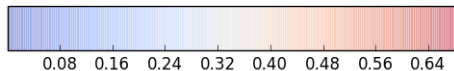
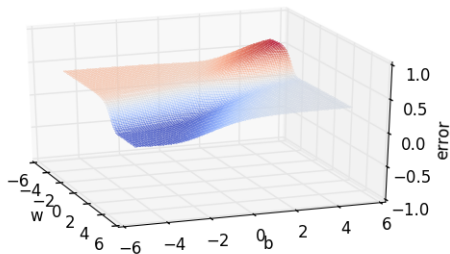
w	b	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

With some guess work and intuition we were able to find the right values for w and b

Let us look at something better than our “guess work” algorithm....

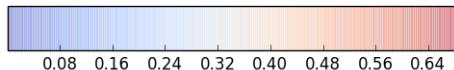
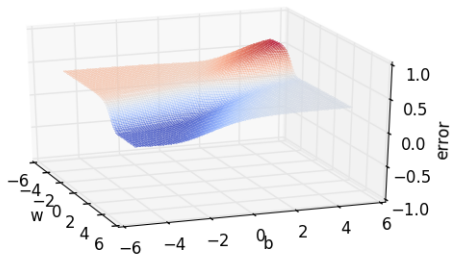
- Since we have only 2 points and 2 parameters (w, b) we can easily plot $\mathcal{L}(w, b)$ for different values of (w, b) and pick the one where $\mathcal{L}(w, b)$ is minimum

Random search on error surface



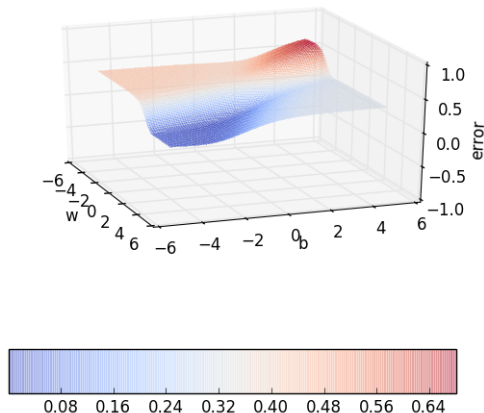
- Since we have only 2 points and 2 parameters (w, b) we can easily plot $\mathcal{L}(w, b)$ for different values of (w, b) and pick the one where $\mathcal{L}(w, b)$ is minimum

Random search on error surface



- Since we have only 2 points and 2 parameters (w, b) we can easily plot $\mathcal{L}(w, b)$ for different values of (w, b) and pick the one where $\mathcal{L}(w, b)$ is minimum
- But of course this becomes intractable once you have many more data points and many more parameters !!

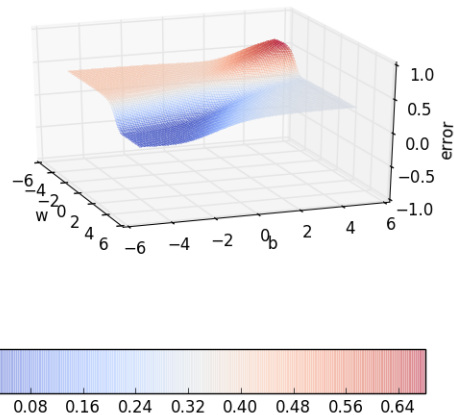
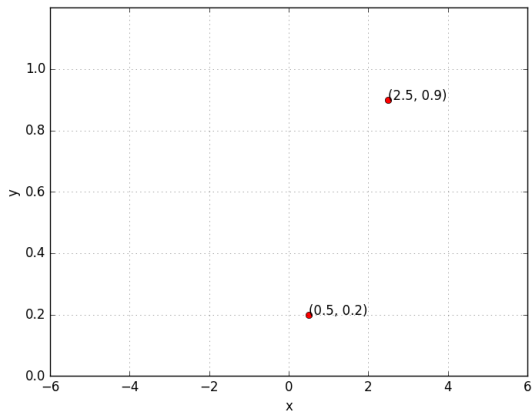
Random search on error surface

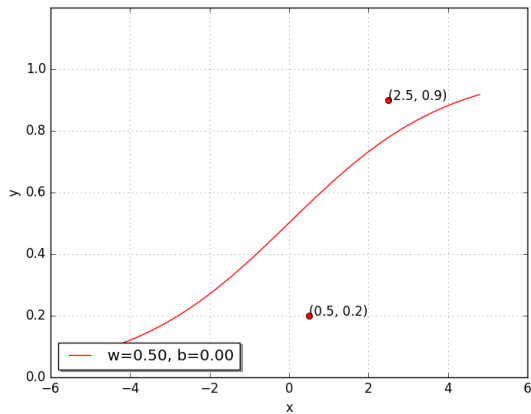


- Since we have only 2 points and 2 parameters (w, b) we can easily plot $\mathcal{L}(w, b)$ for different values of (w, b) and pick the one where $\mathcal{L}(w, b)$ is minimum
- But of course this becomes intractable once you have many more data points and many more parameters !!
- Further, even here we have plotted the error surface only for a small range of (w, b) [from $(-6, 6)$ and not from $(-\infty, \infty)$]

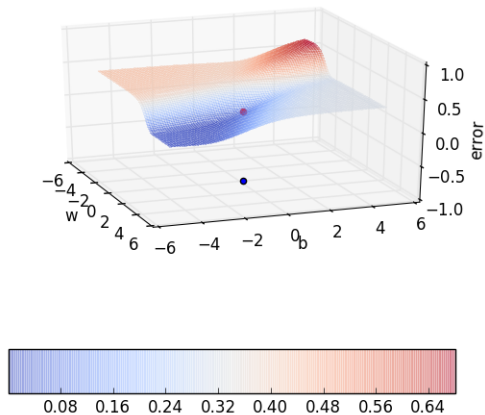
*Let us look at the geometric interpretation of our
“guess work” algorithm in terms of this error surface*

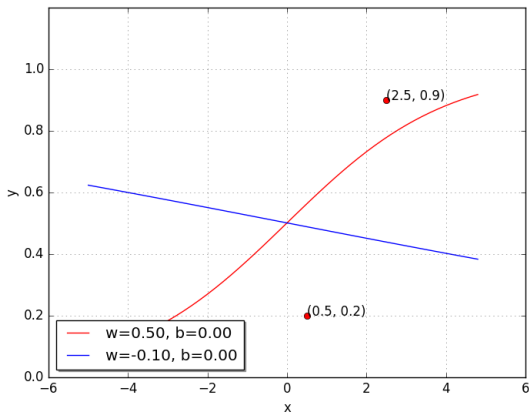
Random search on error surface



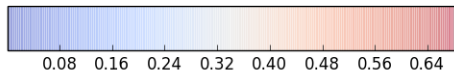
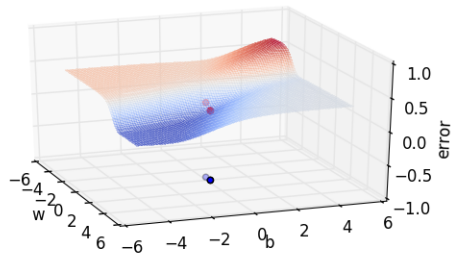


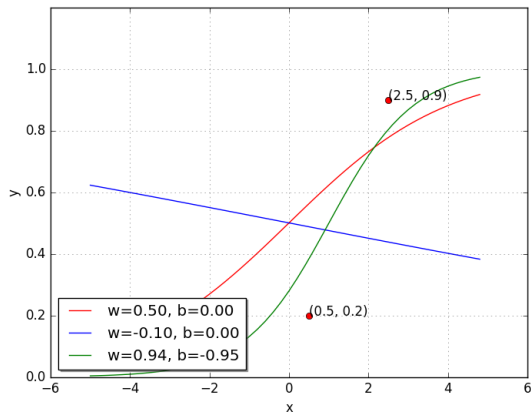
Random search on error surface



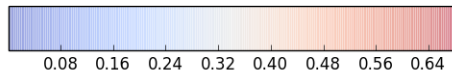
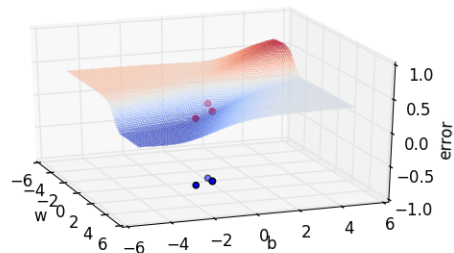


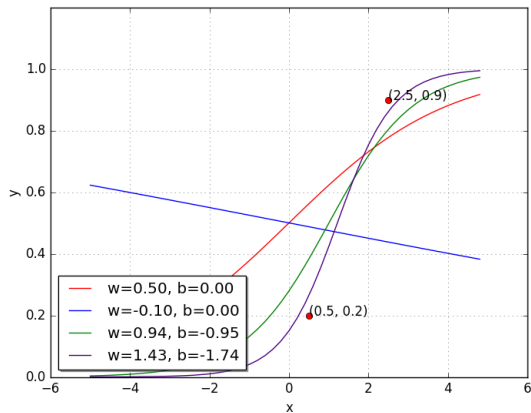
Random search on error surface



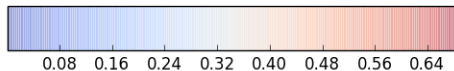
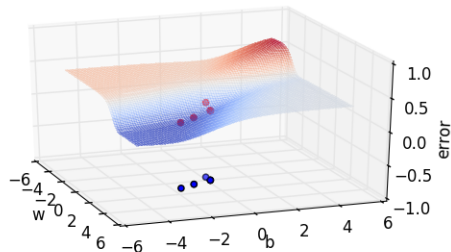


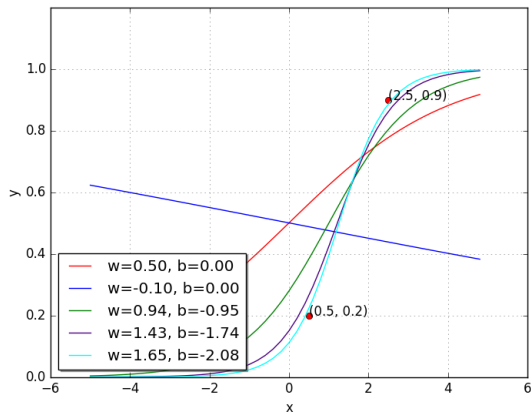
Random search on error surface



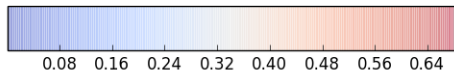
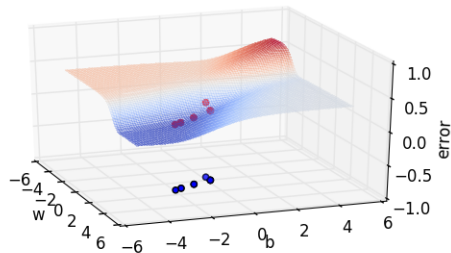


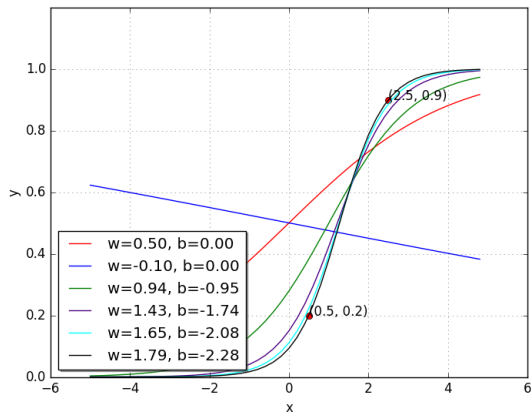
Random search on error surface





Random search on error surface





Random search on error surface

