# Lecture 3 [Logically Weighted & Logistic Regression]

**Linear Regression(recap):**
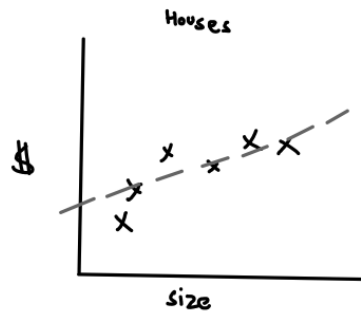
Recap: $(x^{(i)}, y^{(i)})$ — ith example

$$x^{(i)} \in \mathbb{R}^{n+1}, \quad y^{(i)} \in \mathbb{R}, \quad x^{(0)} = 1$$

$m = \text{\# examples}, \quad n = \text{\# features}$

Hypothesis: $h_\theta(x) = \sum\limits_{j=0}^{n} \theta_j x_j = \theta^T x$    "Linear function of features x"

Cost Function: $J(\theta) = \frac{1}{2} \sum\limits_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$   "minimize as function of J to find the parameters $\theta$ for your straight line fit to the data"

Houses

$\$$    x x x x x x x

size

Feature selection algorithms are types of algorithms for automatically deciding what set of features does the best job fitting the data that you have if it's not fit well by a perfectly straight line

**Logically weighted regression(depends on linear regression):**

A way to modify linear regression to make it fit very non-linear functions (so you aren't just fitting straight lines)

In machine learning, we sometimes distinguish between **parametric learning algorithms** and **non-parametric learning algorithms**. In a parametric learning algorithm you fit some fixed set of parameters ( $\theta_i$ ) to data

Locally weighted regression is a non-parametric learning algorithm. The amount of data/parameters you need to keep, grows (linearly in

this case) with the size of the data/training set

With a parametric learning algorithms, no matter how big your training set is, you fit the parameters ( $\Theta_i$ ), then you could erase the training set from your computer memory and make predictions just using the parameters ( $\Theta_i$ )

In a non-parametric learning algorithm, the amount of stuff you need to store around in computer memory or on disk, grows linearly as a function of the training set size

# Locally weighted regression

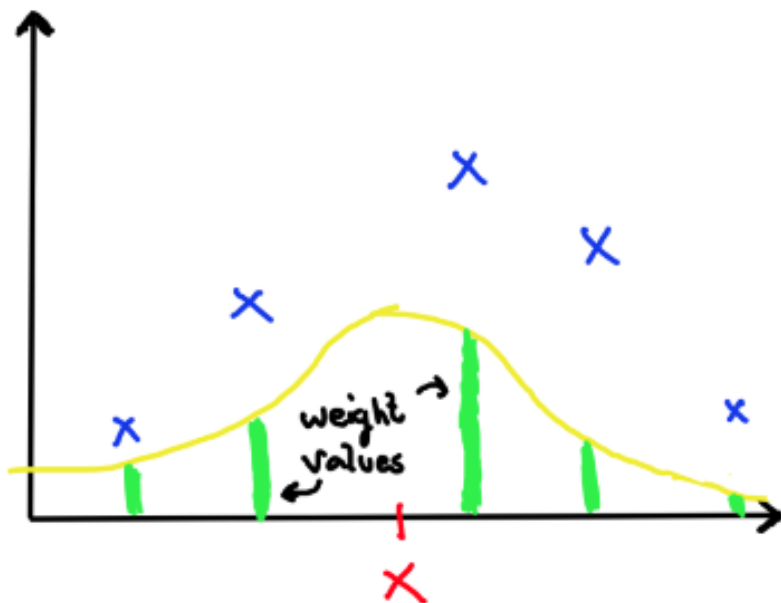Fit $\theta$ to minimize a modified cost function,

$$\sum_{i=1}^{m} w^{(i)} \left( y^{(i)} - \theta^T x^{(i)} \right)^2$$

where $w^{(i)}$ is a "weighting" function

A common/default choice of $w^{(i)}$ is,

$$w^{(i)} = \exp \left( - \frac{(x^{(i)} - x)^2}{\lambda} \right.$$

- If $|x^{(i)} - x|$ is small, $w^{(i)} \approx 1$
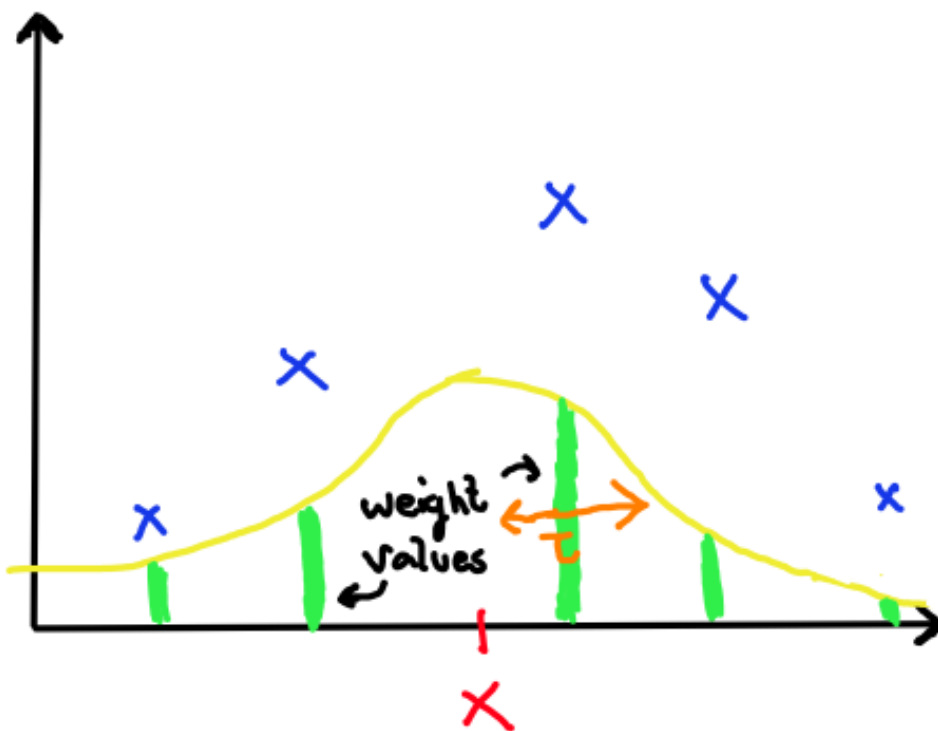- If $|x^{(i)} - x|$ is large, $w^{(i)} \approx 0$



weight → 
values ←

How do you choose the width Gaussian density? How fat/thin should it be? This decides how big a neighborhood should you look in order to decide what's the neighborhood of points you should use to fit this local straight line

A common/default choice of $w^{(i)}$ is,

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

- If $|x^{(i)} - x|$ is small, $w^{(i)} \approx 1$
- If $|x^{(i)} - x|$ is large, $w^{(i)} \approx 0$



"Tau"; $\tau$ - "bandwith parameter"

This is a hyper-parameter of the algorithm

Depending on the choice of tau, you can choose a fatter or a thinner bell-shaped curve which causes you to look in bigger or a narrower window in order to decide how many nearby examples to use in order to fit the straight line

The choice of the bandwith tau has an effect on overfitting and underfitting. If tau is too broad, you end up over-smoothing the data and if tau is too thin, you end up fitting a very jagged fit to the data

Locally weighted linear regression is usually used when you have a relatively low dimensional dataset (when the number of features is not too big) and you have a lot of data and you don't want to think about what features to use

**Probabilistic interpretation (of linear regression):**

Why do we use least squares (squared error)?

Why Least squares?

Assume $\quad y^{(i)} = \Theta^T x^{(i)} + \varepsilon^{(i)}$

$\nwarrow$ "error" : (unmodeled effects, random noise)

$$\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

normal/Gaussian distribution
with mean $0$, variance $\sigma^2$

"is distributed"  "stripped n"  "sigma"

$$P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

$\nearrow$
Probability
Density

- Assume the error terms are I.I.D. (Independently and Identically Distributed)

This implies that :

random variable     mean

$$P(y^{(i)} \mid x^{(i)}; \Theta) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(y^{(i)} - \Theta^T x^{(i)})^2}{2\sigma^2}\right)$$

"parameterized by"

$\nwarrow$ variance

I.e. $\quad y^{(i)} \mid x^{(i)}; \Theta \quad \sim \quad \mathcal{N}(\Theta^T x^{(i)}, \sigma^2)$

If you are willing to make those assumptions then linear regression falls out, almost naturally, of the assumptions you just made, and in particular, under the assumptions you just made:

$$\mathcal{L}(\Theta) = P(\vec{y} \mid x ; \Theta) = \prod_{i=1}^{m} P(y^{(i)} \mid x^{(i)} ; \Theta)$$

"Likelihood of $\Theta$"

"probability of all values of $y$, $(y', \ldots y^m)$, given all $x$'s and parameterized by $\Theta$"

$$= \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(y^{(i)} - \Theta^{T} x^{(i)})^2}{2\sigma^2}\right)$$

What's the difference between Likelihood and Probability?

The Likelihood of the parameters is exactly the same thing as the probability of the data

If you view the value of   P ( y | x ; Θ )   as a function of the parameters holding the data fixed then we call that the Likelihood. If you think of the training set, the data is a fixed thing and then varying parameters Θ, then you're going to use the term Likelihood. Whereas if you view the parameters Θ as fixed and maybe varying the data, you're going to use the term Probability

"Likelihood of the parameters" === "Probability of the data"

"Log Likelihood"

$$\ell(\Theta) = \log \mathcal{L}(\Theta)$$

$$= \log \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(y^{(i)} - \Theta^{T} x^{(i)})^2}{2\,\sigma^2}\right)$$

·Log of a product is equal to the sum of the logs

$$= \sum_{i=1}^{m} \left[ \log \frac{1}{\sqrt{2\pi}\,\sigma} + \log \exp(\cdots) \right]$$

$$= m \log \frac{1}{\sqrt{2\pi}\,\sigma} + \sum_{i=1}^{m} \left[ -\frac{(y^{(i)} - \Theta^{T} x^{(i)})^2}{2\,\sigma^2} \right.$$

# "Maximum Likelihood Estimation" (MLE)

· Choose $\theta$ to maximize $L(\theta)$
  i.e. choose a value of $\theta$ so that
      that value of $\theta$ maximizes the
      probability of the data

To simplify the algebra, rather than maximizing the **L(Θ)** , it's actually easier to maximize the **l(Θ)**
The **log L(Θ)** is a strictly monotonically increasing function, so the value of Θ that maximizes the **l(Θ)**, should be the same as the value of Θ that maximizes the **L(Θ)**, and if you derive the **l(Θ)**

# "Maximum Likelihood Estimation" (MLE)

· Choose $\theta$ to maximize $L(\theta)$

    i.e. choose a value of $\theta$ so that

        that value of $\theta$ maximizes the

        probability of the data

I.e. Choose $\theta$ to minimize $\frac{1}{2}\sum_{i=1}^{m}\left(y^{(i)} - \theta^T x^{(i)}\right)^2 = J(\theta)$

Choosing the value of Θ to minimize the least squares errors is just finding the **MLE** for the parameters Θ under the set of assumptions you made such that the error terms are Gaussian and **IID**

**Logistic regression(classification algorithm, probabilistic interpretation):**

The key steps are:
    1) Make an assumption about P of y given x, parameterized as Θ
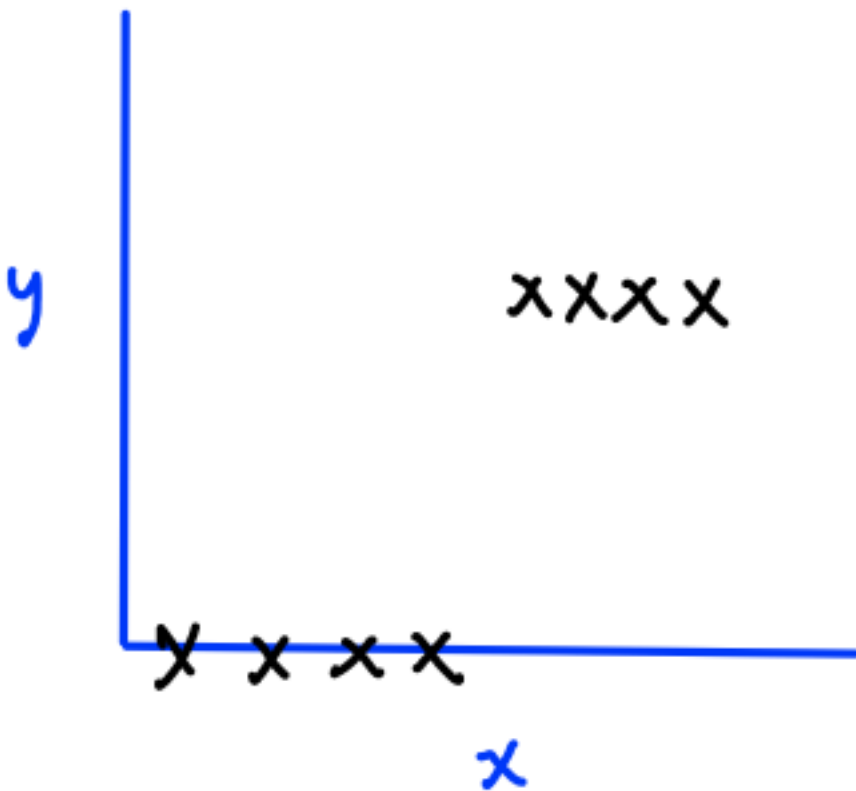    2) Figure out maximum likelihood estimation

Take this framework and apply it to a different type of problem, where the value of y is now either 0 or 1

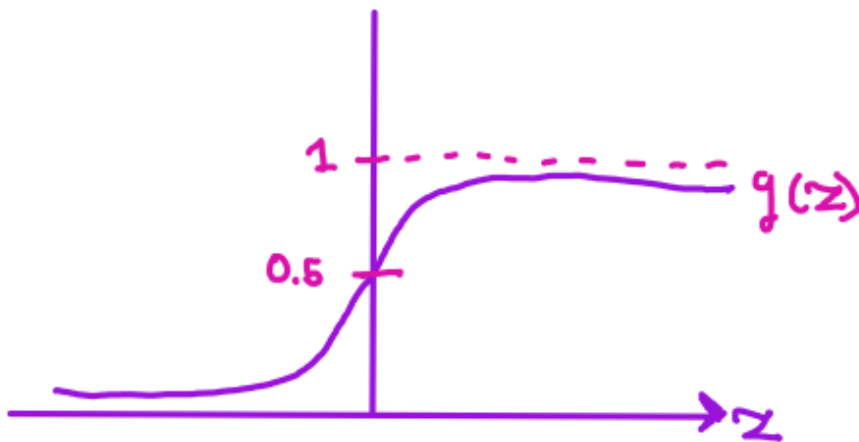# Classification

[Binary Classification]

$$y \in \{0, 1\}$$



Something that's not a good idea is to apply Linear Regression to this data set. You know that for a clasification problem, the values are 0 or 1, so it outputs of negative values or value even greater seems strange

$$\boxed{\text{Logistic Regression}}$$

- You'll want $h_\theta(x) \in [0,1]$

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

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{"sigmoid" or "logistic" function}$$



Let's make some assumptions about the distribution of y given x parameterized by Θ

# Assume the data has the following distribution:

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

Assume what you want your learning algorithm to do is input the features and tell you what's the chance that y is equal to 1

$$\begin{cases} P(y=1 \mid x; \theta) = h_\theta(x) \\ P(y=0 \mid x; \theta) = 1 - h_\theta(x) \end{cases}$$

$$y \in \{0, 1\}$$

Compress

$$P(y \mid x; \theta) = h_\theta(x)^y \left(1 - h_\theta(x)\right)^{1-y}$$

- If $y=1$, $(\text{anything})^0 = 1$, then $P(y \mid x; \theta) = h_\theta(x)$

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

With maximum likelihood estimation, we'll want to find the value of Θ that maximizes the likelihood of the parameters

To make the algebra simpler, take the log of the likelihood and so compute the log likelihood

$$L(\theta) = P(\vec{y} \mid x; \theta)$$

$$= \prod_{i=1}^{m} P(y^{(i)} \mid x^{(i)}; \theta)$$

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

" Log likelihood "

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

$$= \sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log \left(1 - h_\theta(x^{(i)})\right)$$

- Choose the value of $\theta$ to try to maximize $\ell(\theta)$

- The algorithm we're going to use to choose $\theta$ to maximize $\ell(\theta)$, is Gradient ascent / Batch Gradient ascent

update     partial derivative

$$\theta_j := \theta_j + d \frac{d}{d\theta_j} \ell(\theta)$$

The last thing to really flesh out this algorithm is to plug in the definition of the hypothesis. If you do so, you will find that batch gradient ascent is the following:

Batch gradient ascent

$$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} \left( y^{(i)} - h_\theta\left(x^{(i)}\right)\right) x_j^{(i)}$$

$$\underbrace{\qquad\qquad\qquad}_{\frac{\partial}{\partial\theta_j} \ell(\theta)}$$

**Newton's method(algorithm for logistic regression):**

Gradient ascent is a good algorithm, but it takes a lot of iterations(baby steps) for gradient ascent to converge. There's another algorithm called **Newton's method** which allows you to take much larger steps
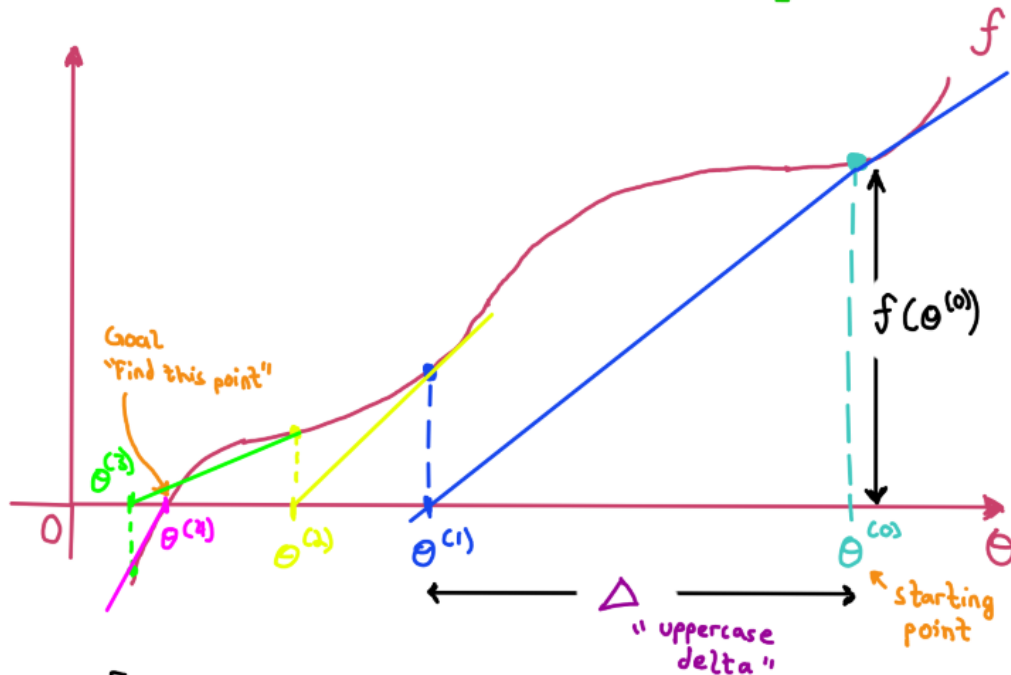
This algorithm is sometimes much faster than gradient ascent for optimizing the value of Θ

[Described with a simplified one-dimensional problem]

where Θ is a real number:

# Newton's method

- Say you have some function $f$, and you want to find a $\theta$ s.t. $f(\theta) = 0$
  This is a problem that Newton's method solves

- What you really want is to maximize $l(\theta)$  (at maximum, first derivative must be 0)
  [i.e. you want to value where the derivative $l'(\theta) = 0$]



Goal
"Find this point"

$f(\theta^{(0)})$

$\leftarrow \quad \Delta \quad \rightarrow$

$\theta^{(0)}$ ← starting point

"uppercase delta"

$\theta^{(3)}$  $\theta^{(2)}$  $\theta^{(2)}$  $\theta^{(1)}$

[this is how 1 iteration of Newton's method will work]

$\delta$ "lowercase delta"

What we'd like to do is solve for the value of delta because one iteration of Newton's method is:

$$\theta^{(1)} := \theta^{(0)} - \Delta$$

$$f'(\theta^{(0)}) = \frac{f(\theta^{(0)})}{\Delta}$$

$$\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\theta^{(t+1)} := \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

Let $f(\theta) = l'(\theta)$  ← we want to find the place where first derivate of $l$ is 0

$$\theta^{(t+1)} := \theta^{(t)} - \frac{l'(\theta^{(t)})}{l''(\theta^{(t)})}$$

Newton's method is a very fast algorithm and enjoys a property called **Quadratic convergence**. What it means is, after one iteration Newton's method has 0.01 error (0.01 away on x-axis from the true value of f = 0), the error could go to 0.0001 error. After 2 iterations it could go to 0.0000001 error

Newton's method, under certain assumptions that functions move not too far from quadratic, the number of significant digits that you have convergedto the minimum, doubles on a single iteration. So this is called **quadratic convergence**

When you get near the minimum, Newton's method converges extremely rapidly, which is why Newton's method requires relatively few iterations

When $\Theta$ is a vector: $(\Theta \in \mathbb{R}^{n+1})$

$$\Theta^{(t+1)} := \Theta^{(t)} + \underbrace{H^{-1}}_{\mathbb{R}^{n+1 \times n+1}} \underbrace{\nabla_\Theta l}_{\text{vector of derivatives } (\mathbb{R}^{n+1})}$$

where H is the Hessian matrix

$$H_{ij} = \frac{\partial^2 l}{\partial \Theta_i \, \partial \Theta_j} \quad \leftarrow \text{defined as the matrix of partial derivatives}$$

The disadvantage of Newton's method is that in **high-dimensional problems**, if $\Theta$ is a vector, then each step of Newton's method is much more expensive because you're either solving linear system equations or having to invert a pretty big matrix

**RULES OF THUMB:**

If the number of parameters is not too big (10, 50, etc) so that the computational cost per iteration is manageable, you would likely use Newton's method because it converges in a very small number of iterations and can be a much faster algorithm than gradient descent, but if you have a very large number of parameters (10,000, 50,000, etc), you would use gradient descent instead