



# Machine learning

## Regression

(adapted from Stanford ML regression notes)

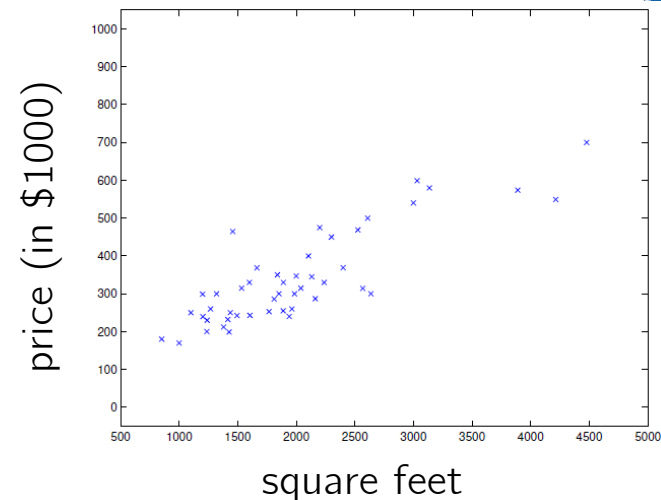
Mohammad-Reza A. Dehaqani

dehaqani@ut.ac.ir

# Prices of houses



Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮



- Given data like this, how can we **learn** to **predict** the prices of other houses, **as a function** of the size of their living areas?
- A pair  $(x^{(i)}, y^{(i)})$  is called a **training example**, and the dataset that we'll be using to learn—a list of  $n$  training examples  $\{(x^{(i)}, y^{(i)}); i=1, \dots, n\}$ —is called a **training set**.
- We used superscript “(i)” in the notation for regression to denote an index into the data set. In other section, we usually use **subscript**.

# Linear Regression

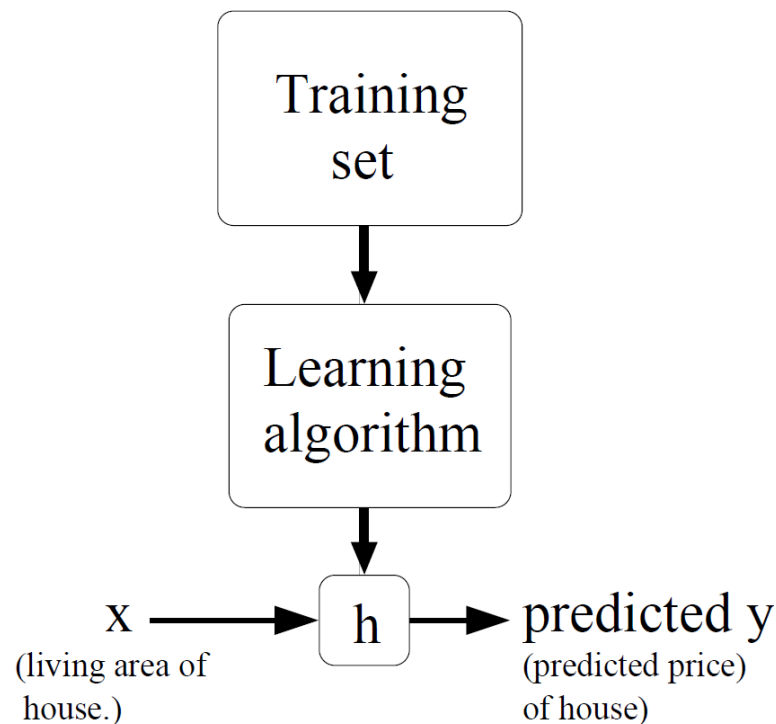


- We approximate  $y$  as a **linear function** of  $x$ :

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

- To perform **supervised learning**, we must decide how we're going to **represent** functions/hypotheses  $h$  in a computer.
- The  $\theta_i$ 's are the **parameters** (also called weights) **parameterizing** the space of linear functions mapping from  $X$  to  $Y$
- Letting  $x_0 = 1$  (this is the **intercept** term), so that the (the new **convention**)

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x$$



$$h : \mathcal{X} \mapsto \mathcal{Y}$$

# How do we pick, or learn, the parameters $\theta$



- One **reasonable** method seems to be to make  $h(x)$  **close to  $y$** , at least for the training examples we have.
- We will define a function that measures, for each value of the  $\theta$ 's, **how close** the  $h(x^{(i)})$ 's are to the corresponding  $y^{(i)}$ 's.
- We define the **cost function** (the ordinary least squares):
$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$
- We want to choose  $\theta$  so as to **minimize**  $J(\theta)$ .

# Gradient descent algorithm to find $\theta$

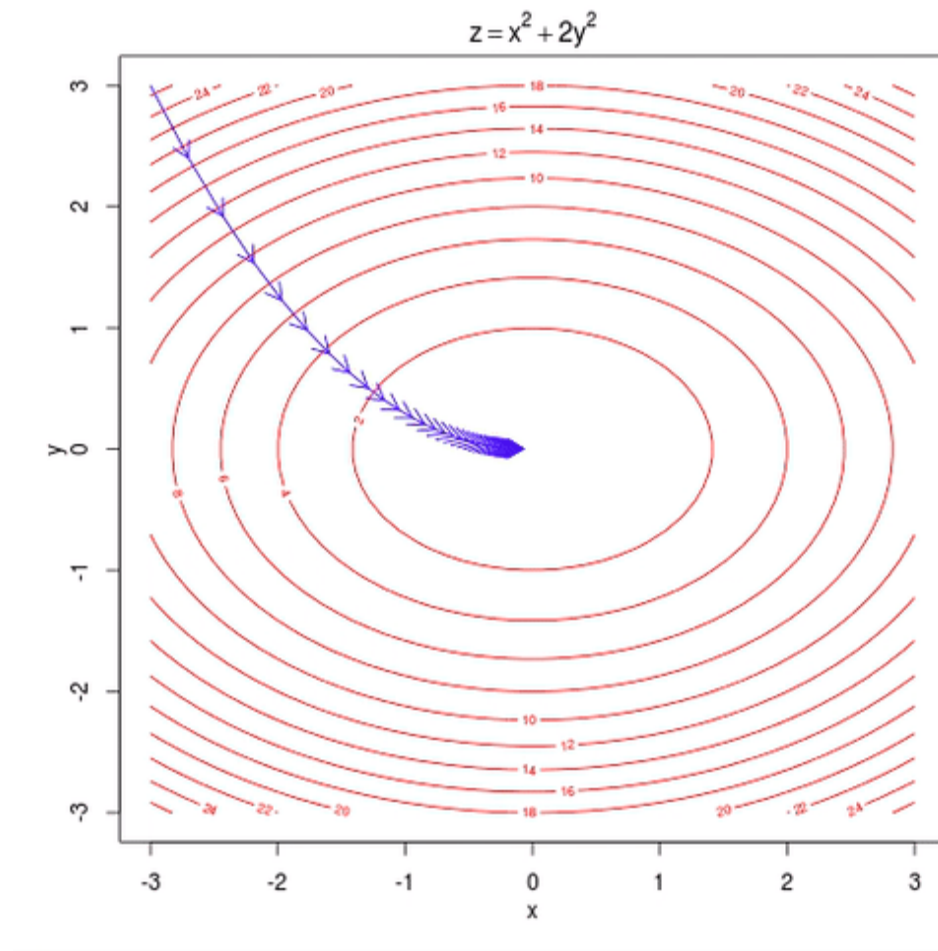


- We update all values of  $\theta_j$ ,  $j = 0, \dots, d$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

- With some “**initial guess**” for  $\theta$ , and that **repeatedly** changes  $\theta$  to make  $J(\theta)$  smaller, until **hopefully** we **converge to a value of  $\theta$**  that minimizes  $J(\theta)$ .
- **$\alpha$**  is called the **learning rate**.
- This is a very **natural algorithm** that repeatedly takes a step in the **direction of steepest decrease** of  $J$

# Gradient descent



# Partial derivative term



- For the case of if we have **only one training example**

$(x, y)$  (neglect the sum in the definition of  $J$ )

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^d \theta_i x_i - y \right) \\ &= (h_{\theta}(x) - y) x_j \end{aligned}$$

- For a single training example, this gives the update rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}.$$

least mean squares (**LMS**)  
update rule or **Widrow-Hoff** learning rule.

# Widrow-Hoff learning rule.



- The magnitude of the update is **proportional** to the **error** term  $(y^{(i)} - h(x^{(i)}))$ ;
- If we are encountering a training example on which our prediction **nearly matches** the actual value of  $y^{(i)}$ , then we find that **there is little need to change the parameters**;
- In contrast, a **larger change** to the parameters will be made if our prediction  $h(x^{(i)})$  has a **large error** (i.e., if it is very far from  $y^{(i)}$ ).





# Modifying method for a training set of more than one example (Solution I)

- **Batch** gradient descent
  - Looks at every example in the **entire training set** on every step, and is called.

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}, \text{ (for every } j \text{)}$$

}

- **Vector notation:**

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)}$$



# Solution II: stochastic gradient descent

- To update the parameters according to the **gradient of the error** with respect to that **single training example only**.

Loop {

for  $i = 1$  to  $n$ , {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j)$$

}

}

$$\theta := \theta + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)}$$

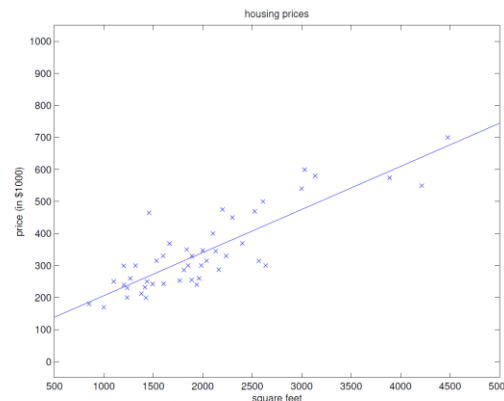
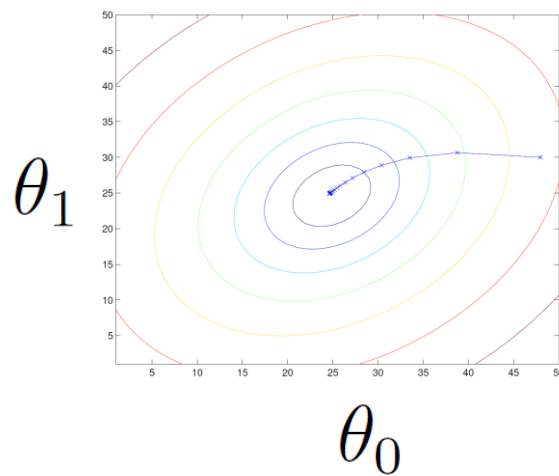
$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)}$$

- Often, stochastic gradient descent gets  $\theta$  “close” to the minimum **much faster** than batch gradient descent.
- When training **set is large**, stochastic gradient descent is **often preferred** over batch gradient descent

# Exmample



- J is a **convex quadratic** function.
- The ellipses shown above are the **contours** of a quadratic function.  $\theta_0 = 71.27, \theta_1 = 0.1345$ .



- Also shown is the **trajectory taken** by gradient descent, which was **initialized at (48,30)**.

# The normal equations



- Performing the minimization explicitly and **without resorting to an iterative algorithm**.
- Define the **design matrix**  $X$  to be the  $n$ -by- $d$  matrix that contains the training examples' input values in its rows
- Also, let  $\vec{y}$  be the  $n$ -dimensional vector containing all the target values from the training set

$$X = \begin{bmatrix} \text{---} & (x^{(1)})^T & \text{---} \\ \text{---} & (x^{(2)})^T & \text{---} \\ & \vdots & \\ \text{---} & (x^{(n)})^T & \text{---} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

# The normal equations



since  $h_{\theta}(x^{(i)}) = (x^{(i)})^T \theta$ ,

$$X\theta - \vec{y} = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(n)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ \vdots \\ h_{\theta}(x^{(n)}) - y^{(n)} \end{bmatrix}$$

$$\frac{1}{2}(X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 = J(\theta)$$



To minimize  $J$ , let's find its derivatives with respect to  $\theta$ .

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\&= \frac{1}{2} \nabla_{\theta} ((X\theta)^T X\theta - (X\theta)^T \vec{y} - \vec{y}^T (X\theta) + \vec{y}^T \vec{y}) \\a^T b &= b^T a \\&= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X) \theta - \vec{y}^T (X\theta) - \vec{y}^T (X\theta)) \\&= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X) \theta - 2(X^T \vec{y})^T \theta) \\\nabla_x b^T x &= b \\&= \frac{1}{2} (2X^T X\theta - 2X^T \vec{y}) \\\nabla_x x^T A x &= 2Ax \\&= X^T X\theta - X^T \vec{y}\end{aligned}$$

# The normal equations



- We set its **derivatives to zero**, and obtain the normal equations:

$$X^T X \theta = X^T \vec{y}$$

- The value of  $\theta$  that minimizes  $J(\theta)$  is given in **closed form** by the equation

$$X \theta = \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

# Probabilistic interpretation



- **Why** linear regression, and the least-squares cost function  $J$ , be a **reasonable choice**?
- Consider following **statistical model**

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

- Where  $\epsilon^{(i)}$  is an **error term** that captures either **unmodeled effects** or **random noise**.
- $\epsilon^{(i)}$  are **distributed IID** (independently and identically distributed) according to a Gaussian distribution

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right).$$

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$



# Maximum likelihood



- The likelihood function:

$$L(\theta) = L(\theta; X, \vec{y}) = p(\vec{y}|X; \theta)$$

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

- The principal of maximum likelihood says that we should **choose  $\theta$**  so as to **make the data as high probability as possible**. I.e., we should choose  $\theta$  to maximize  $L(\theta)$

# Maximize the log likelihood



$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2\end{aligned}$$

- **Maximizing**  $\ell(\theta)$  gives the same answer as **minimizing**

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

- Our final choice of  $\theta$  **did not depend** on what was  $\sigma^2$

# Locally weighted linear regression



- Original linear regression
  1. Fit  $\theta$  to minimize  $\sum_i (y^{(i)} - \theta^T x^{(i)})^2$ .
  2. Output  $\theta^T x$ .
- **Locally weighted** linear regression algorithm
  1. Fit  $\theta$  to minimize  $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$
  2. Output  $\theta^T x$ .
- $w^{(i)}$ 's are **non-negative** valued **weights**.
  - If  $w^{(i)}$  is small, then the  $(y^{(i)} - \theta^T x^{(i)})^2$  error term will be pretty much ignored in the fit.
- A **fairly standard** choice for the weights is

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

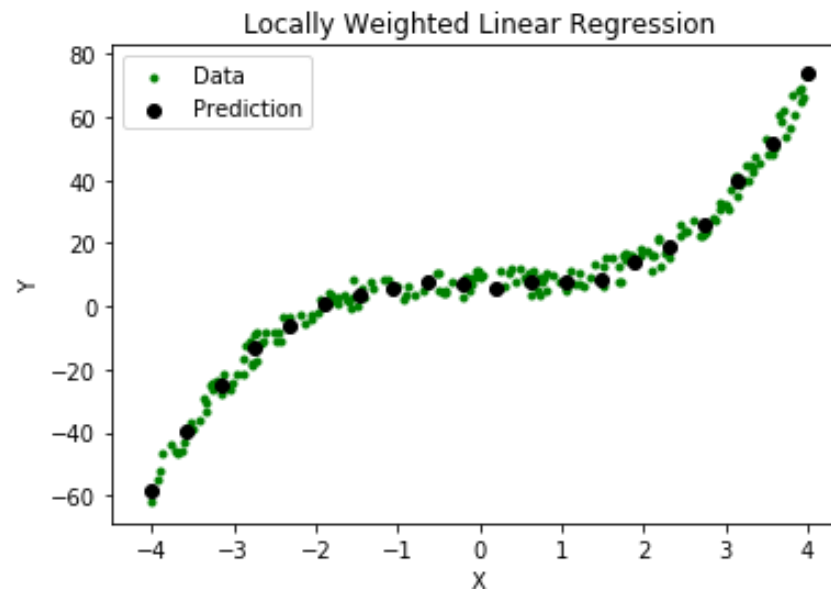
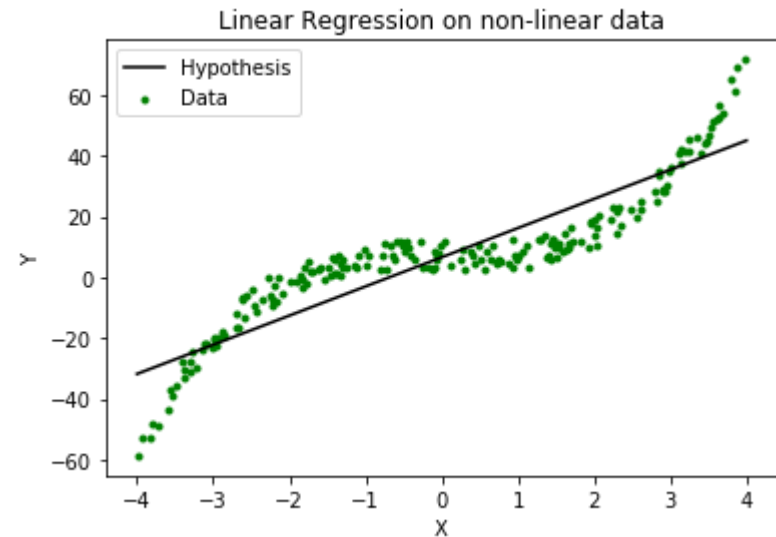
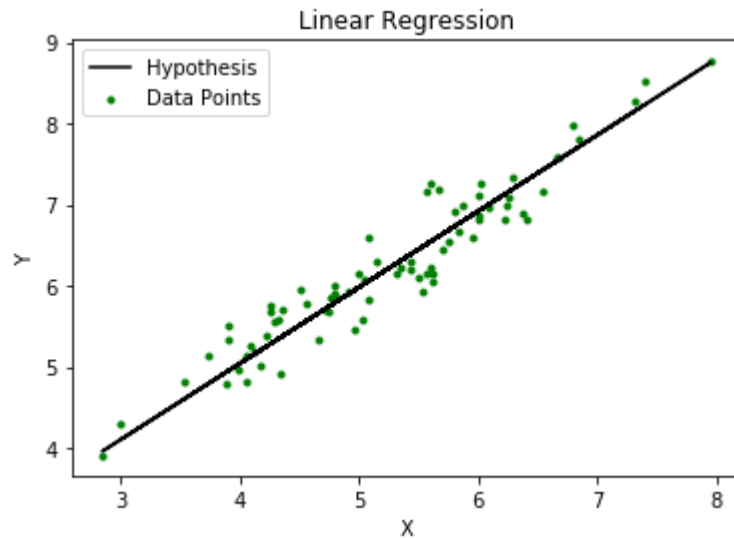
# Non-parametric method

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



- Note that the weights **depend on the particular point  $x$**  at which we're trying to evaluate  $x$ .
- $\theta$  is chosen giving a much higher “weight” to the (errors on) training examples **close to the query** point  $x$ .
- $x$  could be the **position of the center of the peak** in the Bell-shaped function for defining the weights.
- Note that the weights depend on the particular point  $x$  at which we're trying to evaluate  $x$
- $\tau$  is called the **bandwidth parameter**. If  $x$  is vector-valued,  $\tau$  is matrix  $\Sigma$ .
- It is **non-parametric** algorithm:
$$w^{(i)} = \exp(-(x^{(i)} - x)^T \Sigma^{-1} (x^{(i)} - x)/2)$$
- We need **to keep the entire training set** around. The model does not learn a fixed set of parameters as is done in ordinary linear regression
- **Parameters  $\theta$  are computed individually** for each query point

# Example; non-linear relationship between $X$ and $Y$



# Classification and logistic regression



- Let's now talk about the classification problem.
  - This is just like the regression problem, except that the **values  $y$**  we now want to predict take on only a **small number** of **discrete values**.
- For now, we focus on the **binary classification** problem in which  $y$  can take on only two values, 0 and 1.
  - In most cases the binary classifier will also **generalize to the multiple-class** case
- $y^{(i)}$  is called the **label for the training** example.
- Logistic **regression**:
  - We could **approach the classification** problem **ignoring** the fact that  $y$  is discrete-valued

# Logistic function or the sigmoid function.

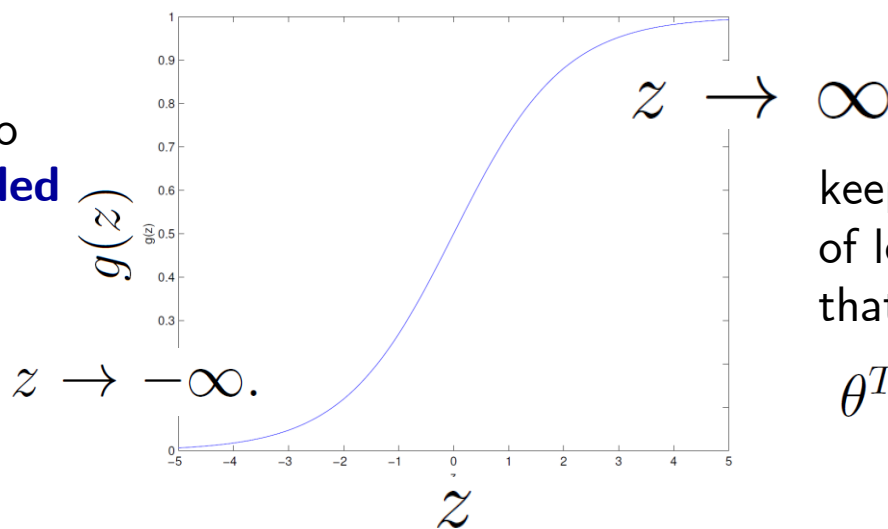


- It also doesn't make sense for  $h_{\theta}(x)$  to take values larger than 1 or smaller than 0 when we know that  $y \in \{0,1\}$ ;
- We will choose:

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

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

$g(z)$ , and hence also  $h(x)$ , is always **bounded** between 0 and 1.



keeping the **convention** of letting  $x_0 = 1$ , so that:

$$\theta^T x = \theta_0 + \sum_{j=1}^d \theta_j x_j$$

# Useful property of the derivative of the sigmoid function,



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

- **Other functions** that smoothly increase from 0 to 1 can also be used
- The choice of the logistic function is a **fairly natural** one: (GLMs, and generative learning algorithms)



# Fitting $\theta$ for logistic regression?



- Setting of **probabilistic assumptions**, and then fit the parameters via **maximum likelihood**.

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

- This can be written more **compactly** as

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

- **IID assumption** on training examples, then write down the **likelihood of the parameters** as

$$\begin{aligned} L(\theta) &= p(\vec{y} \mid X; \theta) \\ &= \prod_{i=1}^n p(y^{(i)} \mid x^{(i)}; \theta) \\ &= \prod_{i=1}^n (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

# Maximizing the log likelihood



$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^n y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))\end{aligned}$$

- We can **use gradient ascent** (Written in vectorial notation)

$$\theta := \theta + \alpha \nabla_{\theta} \ell(\theta)$$

- Start by working with just **one training** example  $(x, y)$ , and take derivatives

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ g'(z) = g(z)(1 - g(z)) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_{\theta}(x)) x_j\end{aligned}$$

# Stochastic gradient ascent rule and perceptron learning



$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

- It is **similar to LMS** update rule; but this is **not** the same algorithm, because  $h_{\theta}(x^{(i)})$  is now defined as a **non-linear function of  $\theta^T x^{(i)}$**
- There is a **deeper reason** on ending up with the same update rule for a rather different algorithm and learning problem. (GLM models)
- Digression: The **perceptron learning** algorithm
  - Modifying the logistic regression method to “**force**” it to output values that are either 0 or 1 or exactly.

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

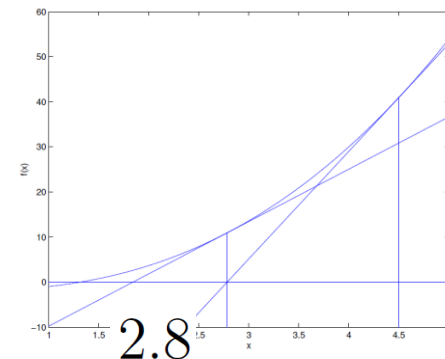
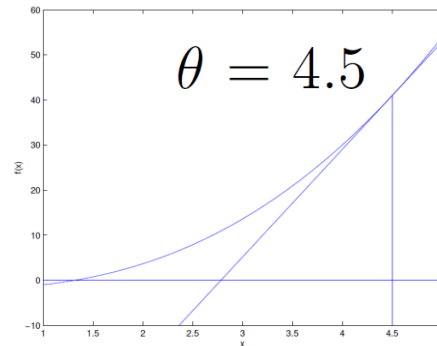
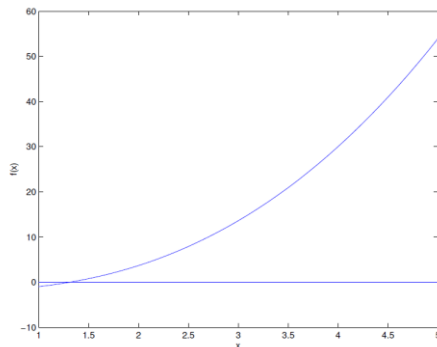
- Using this modified definition of  $g$ , and if we use **the same update rule**, then we have the **perceptron learning algorithm**.

# Newton's algorithm for maximizing $\ell(\theta)$



- **Newton's method** to find a value of  $\theta$  so that  $f(\theta) = 0$ .
- Approximating the function  $f$  via a **linear function** that is tangent to  $f$  at the current guess  $\theta$ ,

$$\theta := \theta - \frac{f(\theta)}{f'(\theta)}$$



- The maxima of  $\ell$  correspond to points where its first derivative  **$\ell'(\theta)$  is zero**. So, by letting  $f(\theta) = \ell'(\theta)$ , we can use the same algorithm to maximize  $\ell$ , and we obtain update rule:

$$\theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

# Newton-Raphson method



- **Vector valued** method in **multidimensional space**

$$\theta := \theta - H^{-1} \nabla_{\theta} \ell(\theta)$$

- $\nabla_{\theta} \ell(\theta)$  is, as usual, the **vector of partial derivatives** of  $\ell(\theta)$  with respect to the  $\theta_i$ 's;
- $H$  is an  $d$ -by- $d$  matrix (actually,  $d+1$ -by- $d+1$ , assuming that we include the intercept term) called the **Hessian**,

$$H_{ij} = \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j}.$$

- It is **faster than gradient descent**; however one step is more **expensive**. Since it requires finding and **inverting an  $d$ -by- $d$  Hessian**
- **Fisher scoring**: Newton's method is applied to maximize the logistic regression

# Generalized Linear Models



- In the regression example, we had

$$y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$$

- In the **classification** one:

$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

- To find some appropriate definitions of  $\mu$  and  $\phi$  **as functions of  $x$  and  $\theta$** .
- Both of these methods are special cases of **a broader family of models**, called Generalized Linear Models

# The exponential family



- A class of distributions

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

- Here,  $\eta$  is called the **natural parameter** (also called the **canonical parameter**) of the distribution;
  - $T(y)$  is the **sufficient statistic** (for the distributions we consider, it will often be the case that  $T(y) = y$ )
  - $a(\eta)$  is the **log partition function**.
  - The quantity  $e^{-a(\eta)}$  essentially plays the role of a **normalization** constant, that makes sure the distribution  $p(y; \eta)$  sums/integrates over  $y$  to 1.
- A fixed choice of  $T$ , **a** and **b** defines a **family** (or set) of distributions that is **parameterized by  $\eta$** ; as we vary  $\eta$ , we then get different distributions within this family.

# Rewrite the Bernoulli distribution



$$\begin{aligned}p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\&= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\&= \exp \left( \left( \log \left( \frac{\phi}{1 - \phi} \right) \right) y + \log(1 - \phi) \right)\end{aligned}$$

- The **natural parameter** is given by  $\eta = \log(\phi/(1-\phi))$ .
- **Interestingly**, if we invert this definition for  $\eta$  by solving for  $\phi$  in terms of  $\eta$ , we obtain  $\phi = 1/(1+e^{-\eta})$ .

- The **familiar sigmoid** function!

$$\begin{aligned}T(y) &= y \\a(\eta) &= -\log(1 - \phi) \\&= \log(1 + e^{\eta}) \\b(y) &= 1\end{aligned}$$

- Bernoulli distribution can be written in the **form of exponential** family, using an appropriate choice of  $T$ ,  $a$  and  $b$



# Rewriting Gaussian distribution.



- Remember that the value of  $\sigma^2$  **had no effect** on our final choice of  $\theta$  and  $h(x)$ . So we choose  $\sigma^2=1$

$$p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right)$$

- We see that the **Gaussian is in the exponential family**,

$$\eta = \mu$$

$$T(y) = y$$

$$a(\eta) = \mu^2/2$$

$$= \eta^2/2$$

$$b(y) = (1/\sqrt{2\pi}) \exp(-y^2/2).$$

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

- There're **many other distributions** that are members of **the exponential family**:
  - The multinomial
  - The Poisson: for modelling **count-data**
  - The gamma and the exponential: for modelling continuous, non-negative random variables, such as **time-intervals**
  - The beta and the Dirichlet: for **distributions over probabilities**
  - and many more.

# Constructing GLMs



- Suppose you would like to build a model to estimate the **number  $y$  of customers arriving in a store** (the **Poisson** distribution usually gives a good model)
  - **Fortunately**, the Poisson is an **exponential family** distribution
  - We would like to predict the value of some random variable  $y$  as a **function of  $x$**
- Three assumption/ design choices
  1.  $y|x; \theta \sim$  **Exponential family( $\eta$ )**. I.e., given  $x$  and  $\theta$ , the distribution of  $y$  follows some exponential family distribution, with **parameter  $\eta$** .
  2. Given  $x$ , our **goal is** to predict the **expected value** of  $T(y)$  **given  $x$** . If  $T(y) = y$ , so this means we would like the prediction  $h_\theta(x)$  output to satisfy  **$h_\theta(x) = E[y|x]$** .
  3. The natural parameter  **$\eta$  and the inputs  $x$**  are related **linearly**:  $\eta = \theta^T x$

# Ordinary Least Squares



- We model the conditional distribution of  $y$  given  $x$  as a Gaussian  $N(\mu, \sigma^2)$  (Here,  $\mu$  may depend  $x$ .)
- We had  $\mu = \eta$

$$\text{Assumption 2} \quad h_{\theta}(x) = E[y|x; \theta]$$

$$y|x; \theta \sim \mathcal{N}(\mu, \sigma^2) = \mu$$

$$\text{Assumption 1} = \eta$$

$$\text{Assumption 3} = \theta^T x.$$

# Logistic Regression in GLM View



- $y \in \{0, 1\}$ . Given that  $y$  is binary-valued, it is natural to use the Bernoulli family of distributions to model the conditional distribution of  $y$  given  $x$

$$\begin{aligned}\eta = \log(\phi/(1 - \phi)) & & h_{\theta}(x) &= E[y|x; \theta] \\ \phi = 1/(1 + e^{-\eta}) & & &= \phi \\ y|x; \theta \sim \text{Bernoulli}(\phi) & & &= 1/(1 + e^{-\eta}) \\ E[y|x; \theta] = \phi & & &= 1/(1 + e^{-\theta^T x})\end{aligned}$$

- This gives us **hypothesis functions of the form**

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

- This one answer to **why we choose the logistic** function form.

# Response and link function.



- The function  $g$  giving the **distribution's mean** as a **function of the natural parameter**  $p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$

$$g(\eta) = \mathbb{E}[T(y); \eta]$$

- The function  $g$  named **canonical response** function.
- Its inverse,  $g^{-1}$ , is called the **canonical link** function

- $\mu = g(\eta); \eta = X\theta$
  - $\eta = X\theta = g^{-1}(\mu);$
- The canonical response function for the **Gaussian** family is just the **identity function**; and the canonical **response** function for the **Bernoulli** is the **logistic function**.

# Softmax Regression



- $y \in \{1, 2, \dots, k\}$ ; thus we will model it as distributed according to a **multinomial distribution**.
- Expressing the multinomial as an exponential family
- $k$  parameters  $\varphi_1, \dots, \varphi_k$  specifying the **probability of each of the outcomes** ( $k-1$  of the  $\varphi_i$ 's uniquely determines the last one)

$$\phi_i = p(y = i; \phi)$$

$$p(y = k; \phi) = 1 - \sum_{i=1}^{k-1} \phi_i$$

- we will **define**  $\mathbf{T}(y) \in \mathbb{R}^{k-1}$

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(3) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, T(k-1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, T(k) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# Multinomial as the exponential family



- Unlike our previous examples, here we **do not have**  $T(y) = y$ ; also,  $T(y)$  is now a  $k-1$  dimensional vector
- $(T(y))_i$  denotes the  $i^{\text{th}}$  element of the vector  $T(y)$ .

$$E[(T(y))_i] = P(y = i) = \phi_i$$

- Multinomial is a member of the **exponential family**.

$$p(y; \phi) = \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \dots \phi_k^{1 - \sum_{i=1}^{k-1} (T(y))_i}$$

$$= \exp((T(y))_1 \log(\phi_1) + (T(y))_2 \log(\phi_2) + \dots + (1 - \sum_{i=1}^{k-1} (T(y))_i) \log(\phi_k))$$

$$= \exp((T(y))_1 \log(\phi_1/\phi_k) + (T(y))_2 \log(\phi_2/\phi_k) + \dots + (T(y))_{k-1} \log(\phi_{k-1}/\phi_k) + \log(\phi_k))$$

$$b(y) = 1 \quad a(\eta) = -\log(\phi_k) \quad \eta = \begin{bmatrix} \log(\phi_1/\phi_k) \\ \log(\phi_2/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix}$$

$$= \underbrace{b(y)} \exp(\underbrace{\eta^T}_{\text{row vector}} \underbrace{T(y)}_{\text{column vector}} - \underbrace{a(\eta)})$$

# softmax function.



link function is  $\eta_i = \log \frac{\phi_i}{\phi_k}$ . for  $i = 1, \dots, k$ ,  
 $\eta_k = \log(\phi_k/\phi_k) = 0$ .

- To invert the link function and derive the **response function**,

$$e^{\eta_i} = \frac{\phi_i}{\phi_k}$$

$$\phi_k e^{\eta_i} = \phi_i$$

$$\phi_k \sum_{i=1}^k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1$$

$$\phi_k = 1 / \sum_{i=1}^k e^{\eta_i} \xrightarrow{\phi_k e^{\eta_i} = \phi_i} \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$



# Using Assumption 3: the $\eta_i$ 's are linearly related to the $x$ 's.



- We can also define  $\theta_k = 0$ , so that  $\eta_k = \theta_k^T x = 0$

$$p(y = i|x; \theta) = \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}$$

- Where  $y \in \{1, \dots, k\}$ , is called **softmax regression**. It is a **generalization** of logistic regression.

$$h_{\theta}(x) = E[T(y)|x; \theta] = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix} = \begin{bmatrix} \frac{\exp(\theta_1^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \frac{\exp(\theta_2^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \vdots \\ \frac{\exp(\theta_{k-1}^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \end{bmatrix}$$

- Hypothesis will **output the estimated probability** that  $p(y = i|x; \theta)$ , for every value of  $i = 1, \dots, k$

# log-likelihood



$$p(y; \phi) = \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \dots \phi_k^{1 - \sum_{i=1}^{k-1} (T(y))_i}$$
$$\ell(\theta) = \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \theta) \quad \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

$$= \sum_{i=1}^n \log \prod_{l=1}^k \left( \frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{(T(y^{(i)}))_l}$$

- We can now obtain the maximum likelihood estimate of the parameters by maximizing  $\ell(\theta)$  in terms of  $\theta$ , using a **method such as gradient ascent or Newton's** method.