

Learning From Data

Lecture 2: Linear Regression & Logistic Regression

Yang Li yangli@sz.tsinghua.edu.cn

September 24, 2021

Today's Lecture

Supervised Learning (Part I)

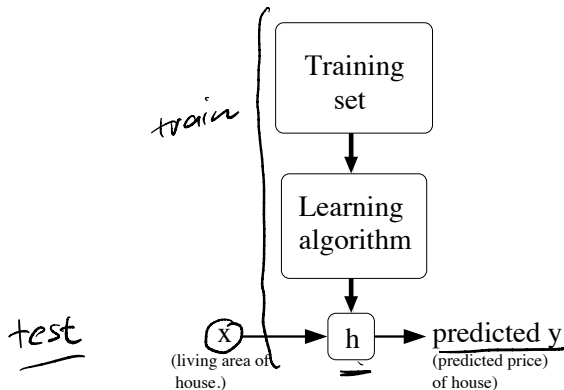
- ▶ Linear Regression
- ▶ Binary Classification
- ▶ Multi-Class Classification

Review: Supervised Learning

- ▶ Input space: \mathcal{X} , Target space: \mathcal{Y}
 \mathbb{R}^d

Review: Supervised Learning

- ▶ Input space: \mathcal{X} , Target space: \mathcal{Y} (x, y)
- ▶ Given training examples, we want to learn a **hypothesis** function $h: \mathcal{X} \rightarrow \mathcal{Y}$ so that $h(x)$ is a "good" predictor for the corresponding y .



Review: Supervised Learning

\mathcal{Y} - target space

- ▶ y is discrete (categorical): **classification problem**
- ▶ y is continuous (real value): **regression problem**

Linear Regression

Linear Regression Model

Ordinary Least Square

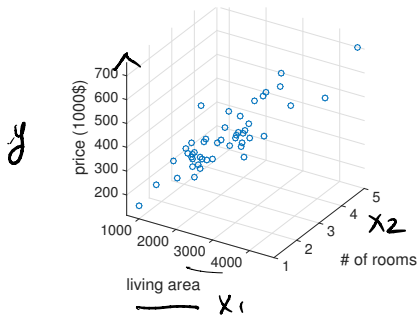
Maximum Likelihood Estimation

y is continuous

Linear Regression

Example: predict Portland housing price

<u>Living area (ft^2)</u>	# bedrooms	<u>Price (\$1000)</u>
x_1	x_2	y
2104	3	400
1600	3	330
2400	3	369
\vdots	\vdots	\vdots



Linear Approximation

$$\begin{array}{l} \text{input space} \quad \text{output space} \\ \mathcal{X} = \mathbb{R}^2 \quad \mathcal{Y} = \mathbb{R} \quad h: \mathcal{X} \rightarrow \mathcal{Y} \\ \hline \mathcal{Y} \in \mathcal{Y} \end{array}$$

A linear model

$$h(x) = \underline{\theta}_0 + \underline{\theta}_1 x_1 + \underline{\theta}_2 x_2$$

θ_i 's are called **parameters**.

Linear Approximation

A linear model

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

θ_i 's are called **parameters**.

Using vector notation,

$$\underline{h(x)} = \underline{\theta}^T \underline{x}, \quad \text{where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

Alternative Notation

$$h(x) = \underline{w_1}x_1 + \underline{w_2}x_2 + \underline{b}$$

w_1, w_2 are called weights, b is called the bias

$$h(x) = w^T x + b, \quad \text{where } w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

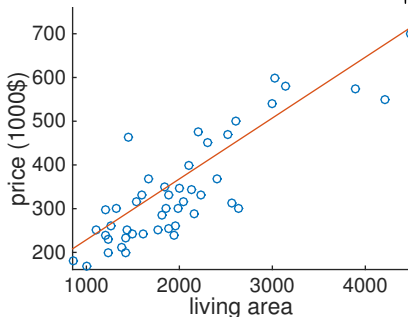
Apply model to new data

Suppose we have the optimal parameters $\underline{\theta}$, e.g.

```
> h = LinearRegression().fit(X, y)
> theta = h.coef
array([89.60, 0.1392, -8.738])
      theta_0 theta_1 theta_2
```

make a prediction of new feature x :

$$\hat{y} = h_{\theta}(x) = \theta^T x \quad \text{evaluating the model} \\ \rightarrow \text{prediction}$$



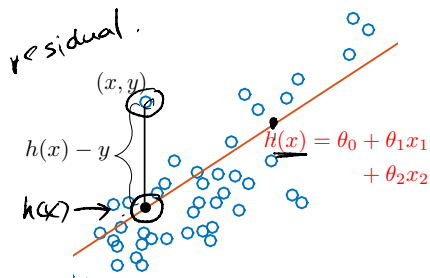
Model Estimation

How to estimate model parameters θ (or \underline{w} and \underline{b}) from data?

Model Estimation

How to estimate model parameters θ (or w and b) from data?

Least Square Estimation

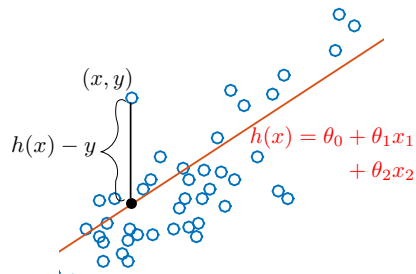


geometric approach

Model Estimation

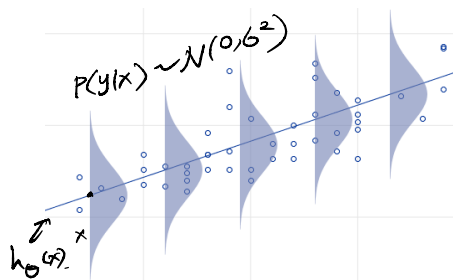
How to estimate model parameters θ (or w and b) from data?

Least Square Estimation



geometric approach

Maximum Likelihood Estimation



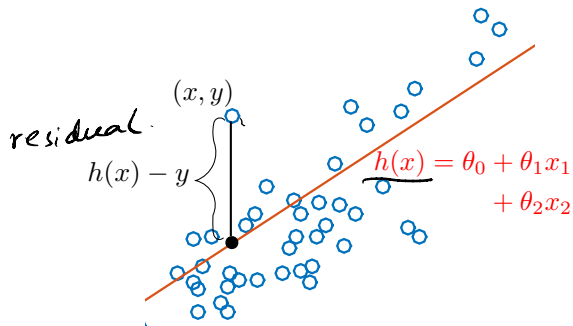
Probabilistic approach

Ordinary Least Square

Cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

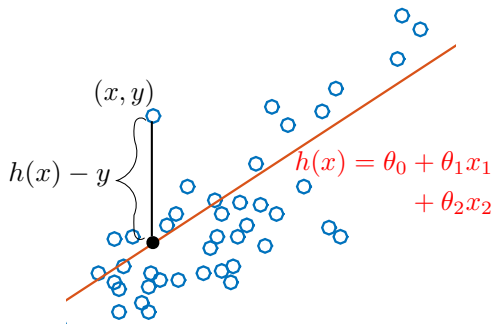
of training samples
prediction
true label



Ordinary Least Square

Cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$



Ordinary Least Square

Cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

Ordinary Least Square

Cost function:

$$\underline{J(\theta)} = \frac{1}{2} \sum_{i=1}^m \underbrace{(h(x^{(i)}) - y^{(i)})^2}_{\theta^T x^{(i)}}$$

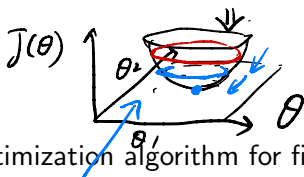
The **ordinary Least square problem** is:

$$\begin{aligned} & \min_{\theta} \underbrace{J(\theta)} \\ & = \min_{\theta} \frac{1}{2} \sum_{i=1}^m \underbrace{(h(x^{(i)}) - y^{(i)})^2} \end{aligned}$$

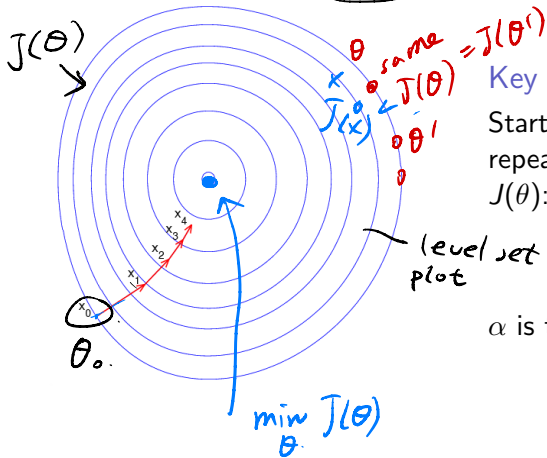
How to minimize $J(\theta)$?

- ▶ Numerical solution: gradient descent, Newton's method
- ▶ Analytical solution: normal equation

Gradient descent



A first-order iterative optimization algorithm for finding the minimum of a function $J(\theta)$.



Key idea

Start at an initial guess, repeatedly change θ to decrease $J(\theta)$:

$$\theta := \theta - \underbrace{\alpha \nabla J(\theta)}_{\text{gradient}}$$

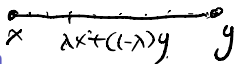
α is the **learning rate**

0.1

Review: Convex function

Vector space S , a subset $C \subseteq S$ is convex if for any $x, y \in C$, its $\overline{xy} \subseteq C \Rightarrow$ for any $0 \leq \lambda \leq 1$, affine combination $\lambda x + (1-\lambda)y \in C$.

Definition



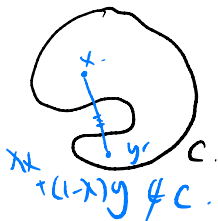
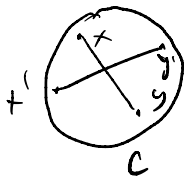
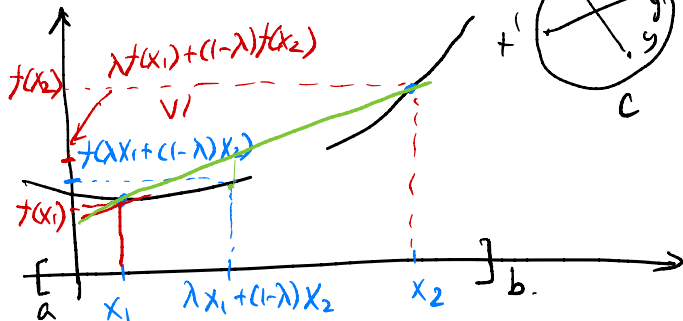
A function $f(x)$ is **convex** on a convex set C if for any $x_1, x_2 \in C$ and $0 \leq \lambda \leq 1$,

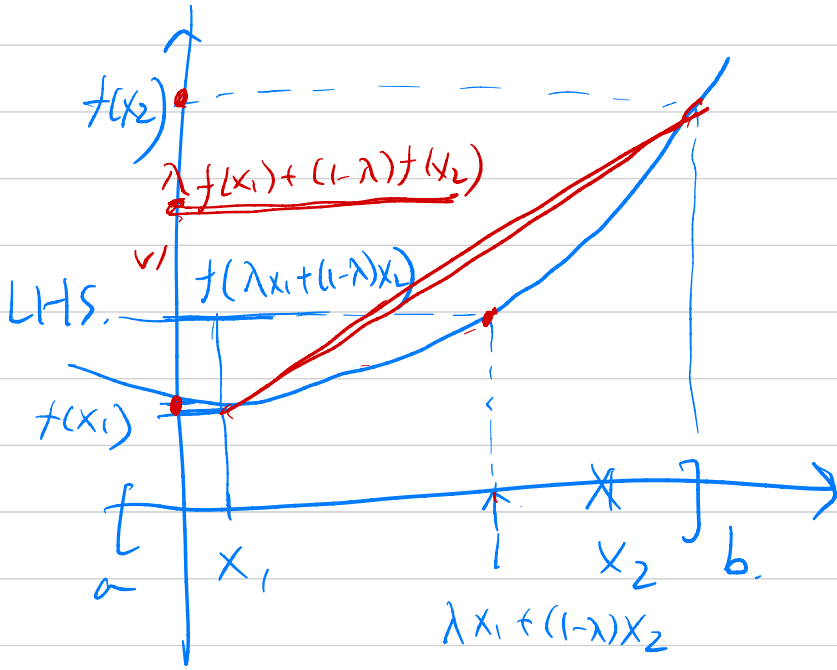
$$f: C \rightarrow \mathbb{R}$$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

e.g. C is an interval $[a, b]$

$\lambda = 0.6$





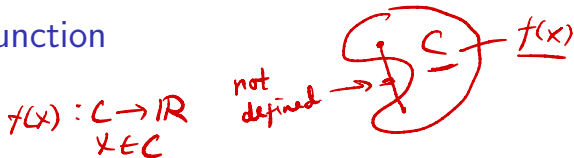
Review: Convex function

Definition

A function $f(x)$ is **convex** on a convex set C if for any $x_1, x_2 \in C$ and $0 \leq \lambda \leq 1$,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

e.g. C is an interval $[a, b]$



Theorem

If $J(\theta)$ is convex, gradient descent finds the global minimum.

For the ordinary least square problem,

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

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

$\begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \leftarrow \# \text{ features.}$

, where $\frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$

$$= \frac{1}{2} \sum_{i=1}^m 2 (\theta^T x^{(i)} - y^{(i)}) \frac{\partial (\theta^T x^{(i)} - y^{(i)})}{\partial \theta_j}$$

$$= \sum_{i=1}^m x_j^{(i)} (\theta^T x^{(i)} - y^{(i)})$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \sum_{k=1}^n \theta_k x_k^{(i)} &= \sum_{k=1}^n \frac{\partial}{\partial \theta_j} \theta_k x_k^{(i)} \\ &= \frac{\partial}{\partial \theta_j} \theta_j x_j^{(i)} \\ &= x_j^{(i)} \end{aligned}$$

For the ordinary least square problem,

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

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}, \text{ where } \frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \right]$$
$$= \underbrace{\sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}}_{}$$

Gradient descent for ordinary least square

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$

Gradient of cost function: $\nabla J(\theta)_j = \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}$

Gradient descent update: $\theta_j \leftarrow \theta_j - \alpha \nabla J(\theta)_j$

Batch Gradient Descent

Repeat until convergence { ←

$\theta_j = \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$ for every j ←

while not converged:

for j in range(n):

$\theta_j = \dots$

Gradient descent for ordinary least square

Gradient of cost function: $\nabla J(\theta)_j = \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}$

Gradient descent update: $\theta := \theta - \alpha \nabla J(\theta)$

Batch Gradient Descent

Repeat until convergence {

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

}

$O(m)$

θ is only updated after we have seen all m training samples.

Batch gradient descent

```
Repeat until convergence{  
   $\theta_j = \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))x_j^{(i)}$  for every j  
}
```

for every j

for every $i=1, \dots, m$
 $\theta_j += \alpha (y^{(i)} - h_{\theta}) x_j^{(i)}$

θ_j .

Stochastic gradient descent

```
Repeat until convergence{  
  for  $i=1 \dots m$  {  
     $\theta_j = \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)}))x_j^{(i)}$  for every j  
  }  
}
```

$\theta = [\theta_1, \dots, \theta_n]$.

θ is updated each time a training example is read

Batch gradient descent

```
Repeat until convergence{  
   $\theta_j = \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))x_j^{(i)}$  for every j  
}
```

Stochastic gradient descent

```
Repeat until convergence{  
  for  $i = 1 \dots m$  {  
     $\theta_j = \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)}))x_j^{(i)}$  for every j  
  }  
}
```

mini-batch



θ is updated each time a training example is read

- ▶ Stochastic gradient descent gets θ close to minimum much faster
- ▶ Good for regression on large data

Minimize $J(\theta)$ Analytically

The matrix notation

$$\underline{X} = \begin{matrix} \left. \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \vdots & & \\ \text{---} & \text{---} & \text{---} \end{matrix} \right\} \begin{matrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{matrix} \end{matrix}, \quad \underline{\vec{y}} = \begin{matrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{matrix}$$

Handwritten annotations: $x^{(i)}$ above the matrix columns, n under the columns, m to the left of the rows, and $(m \times 1)$ to the right of the vector \vec{y} .

X is called the **design matrix**.

Minimize $J(\theta)$ Analytically

$$z = X\theta - y = \begin{bmatrix} \theta^T x^{(1)} \\ \theta^T x^{(2)} \\ \theta^T x^{(3)} \\ \vdots \\ \theta^T x^{(m)} \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} \theta^T x^{(1)} - y^{(1)} \\ \vdots \\ \theta^T x^{(i)} - y^{(i)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{bmatrix}$$

The matrix notation

$$z^T z = \sum_{i=1}^m z_i^2$$

$z \in \mathbb{R}^m$

$$\frac{1}{2} (X\theta - y)^T (X\theta - y) \quad X = \begin{bmatrix} - (x^{(1)})^T & - \\ - (x^{(2)})^T & - \\ \vdots & \vdots \\ - (x^{(m)})^T & - \end{bmatrix} \begin{bmatrix} \theta \end{bmatrix}, \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

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

X is called the **design matrix**. The least square function can be written as

$$J(\theta) = \frac{1}{2} \underline{(X\theta - y)^T (X\theta - y)}$$

Compute the gradient of $J(\theta) : \frac{1}{2}(Qx-y)^T(Qx-y)^T$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \left[\frac{1}{2}(X\theta - y)^T(X\theta - y) \right]$$

Compute the gradient of $J(\theta)$:

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \left[\frac{1}{2} (\underline{X\theta} - y)^T (X\theta - y) \right]$$

Hint: (let $\underline{x} \leftarrow \theta$, $\underline{Q} \leftarrow X$)

$x \in \mathbb{R}^{m+n}$
 $x, y \in \mathbb{R}^n$, $Q \in \mathbb{R}^{n \times n}$. $l: \mathbb{R}^n \rightarrow \mathbb{R}$.

$$\begin{aligned} l(x) &= \frac{1}{2} (\underline{Qx - y})^T (Qx - y) \\ &= \frac{1}{2} (x^T Q^T - y^T) (Qx - y) \\ &= \frac{1}{2} (x^T Q^T Q x - \underline{y^T Q x} - \underline{x^T Q^T y} + y^T y) \\ &= \frac{1}{2} x^T Q^T Q x - y^T Q \cdot x + \frac{1}{2} y^T y \\ \nabla_x l(x) &= \frac{1}{2} \cdot \frac{\partial (x^T Q^T Q x)}{\partial x} - \frac{\partial (y^T Q \cdot x)}{\partial x} + 0. \\ &= \frac{1}{2} \cdot (Q^T Q + Q^T Q^T) \cdot x - Q^T y \\ &= \underline{Q^T Q \cdot x - Q^T y} \end{aligned}$$

Compute the gradient of $J(\theta)$:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \left[\frac{1}{2} (X\theta - y)^T (X\theta - y) \right] \\ &= \end{aligned}$$

Compute the gradient of $J(\theta)$:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \left[\frac{1}{2} (X\theta - y)^T (X\theta - y) \right] \\ &= \underline{X^T X \theta - X^T y}\end{aligned}$$

quadratic form
 $\sum_{i=1}^m (\theta^T x^i - y^i)^2$

Since $J(\theta)$ is convex, x is a global minimum of $J(\theta)$ when $\nabla J(\theta) = 0$.

Compute the gradient of $J(\theta)$:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \left[\frac{1}{2} (X\theta - y)^T (X\theta - y) \right] \\ &= X^T X \theta - X^T y = 0 \\ &\quad \theta = (X^T X)^{-1} X^T y\end{aligned}$$

Since $J(\theta)$ is **convex**, x is a global minimum of $J(\theta)$ when

$$\nabla J(\theta) = 0.$$

The Normal equation

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

Compute the gradient of $J(\theta)$:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \left[\frac{1}{2} (X\theta - y)^T (X\theta - y) \right] \\ &= X^T X\theta - X^T y\end{aligned}$$

Since $J(\theta)$ is **convex**, x is a global minimum of $J(\theta)$ when $\nabla J(\theta) = 0$.

The Normal equation

$$\begin{aligned}X \in \mathbb{R}^{m \times n} &\rightarrow \theta = X^{-1} y. \\ X \in \mathbb{R}^{m \times n} &\rightarrow \underbrace{(X^T X)^{-1} X^T}_{\theta} y. \\ \theta &= \underbrace{(X^T X)^{-1} X^T}_y y\end{aligned}$$

$(X^T X)^{-1} X^T$ is called the **Moore-Penrose pseudoinverse** of X

Which method to use?

gradient descent	normal equation
iterative solution	exact solution

Which method to use?

gradient descent	<u>normal equation</u>
iterative solution	exact solution
need to choose proper learning parameter α for cost function to converge <i>steepest GD.</i>	

Which method to use?

$X: (m \times n)$ $n \times n$
 $(X^T X)^{-1}$
Avoid numerical issue: add small perturbation:
 $(X^T X + \lambda I)^{-1}$
 $\begin{bmatrix} \lambda & & \\ & \lambda & \\ & & \lambda \end{bmatrix}$

gradient descent	normal equation
iterative solution	exact solution
need to choose proper learning parameter α for cost function to converge	<u>numerically unstable</u> when <u>X</u> is <u>ill-conditioned</u> . e.g. features are highly correlated

Which method to use?

gradient descent	normal equation
iterative solution	exact solution
need to choose proper learning parameter α for cost function to converge	numerically unstable when X is ill-conditioned. e.g. features are highly correlated
works well for large number of samples m	

Which method to use?

gradient descent	normal equation
iterative solution	exact solution
need to choose proper learning parameter α for cost function to converge	numerically unstable when X is ill-conditioned. e.g. features are highly correlated
works well for large number of samples m	solving equation is slow when <u>m is large</u>

Minimize $J(\theta)$ using Newton's Method

$f(\theta)$,

Numerically solve for θ in $\nabla_{\theta} J(\theta) = 0$

Newton's method

Solves real functions $f(x) = 0$ by iterative approximation:

- ▶ Start an initial guess x
- ▶ Update x until convergence

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

Minimize $J(\theta)$ using Newton's Method

Geometric intuition of Newton's method

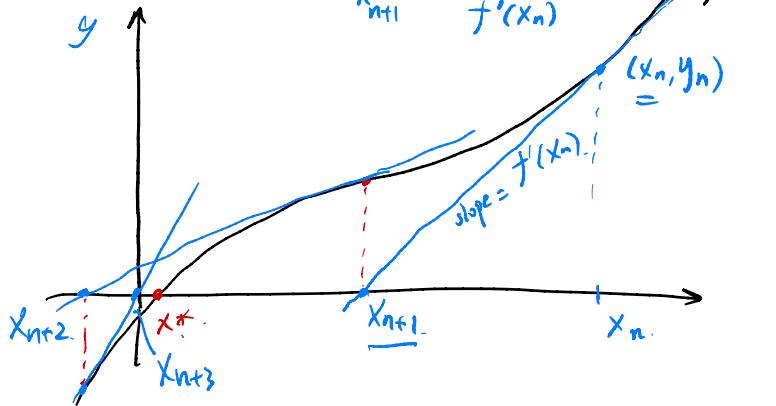
At step $n + 1$:

- ▶ Find tangent line of f at (x_n, y_n)
- ▶ $x_{n+1} \leftarrow$ x-intercept of the tangent line
- ▶ $y_{n+1} \leftarrow f(x_{n+1})$

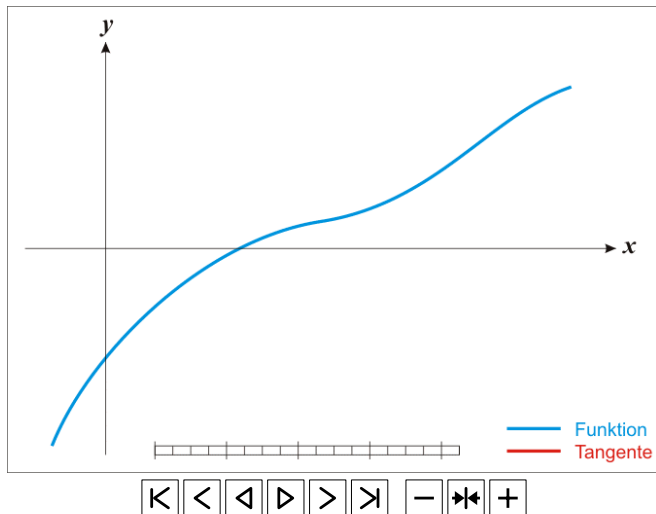
$x_n \leftarrow x$ at n th iteration

$$y = f'(x)(x - x_n) + f(x_n)$$
$$x = -\frac{f(x_n)}{f'(x_n)} + x_n$$

$$x_{n+1} \leftarrow x_n - \frac{f(x_n)}{f'(x_n)}$$



Newton's Method Demo



https://en.wikipedia.org/wiki/File:NewtonIteration_Ani.gif

Minimize $J(\theta)$ using Newton's Method

Newton's method for optimization $\min_{\theta} J(\theta)$

Use Newton's method to solve $\nabla_{\theta} J(\theta) = 0$:

- ▶ θ is one-dimensional:

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

Handwritten annotations: $\frac{\partial}{\partial \theta} J(\theta)$ with an arrow pointing to $J'(\theta)$, and $\frac{\partial^2}{\partial \theta^2} J(\theta)$ with an arrow pointing to $J''(\theta)$.

Minimize $J(\theta)$ using Newton's Method

Newton's method for optimization $\min_{\theta} J(\theta)$

Use Newton's method to solve $\nabla_{\theta} J(\theta) = 0$:

- ▶ θ is one-dimensional:

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

- ▶ θ is n -dimensional:

$$\theta = \theta - \underline{H^{-1}(\theta)} \underline{\nabla J(\theta)}$$

where H is the Hessian matrix of $J(\theta)$.

a.k.a Newton-Raphson method

Hessian matrix,

$H(\theta)$

$$= \begin{bmatrix} \frac{\partial^2 J(\theta)}{\partial \theta_1^2} & \frac{\partial^2 J(\theta)}{\partial \theta_1 \partial \theta_2} & \dots \\ \frac{\partial^2 J(\theta)}{\partial \theta_2 \partial \theta_1} & \dots & \dots \\ \vdots & \dots & \dots \\ \frac{\partial^2 J(\theta)}{\partial \theta_n \partial \theta_1} & \dots & \frac{\partial^2 J(\theta)}{\partial \theta_n^2} \end{bmatrix}$$

Newton's Method for Optimization

```
Initialize  $\theta$   
While  $\theta$  has not coveredged {  
   $\theta := \theta - \underline{H^{-1}(\theta)} \nabla J(\theta)$   
}
```

Newton's Method for Optimization

```
Initialize  $\theta$ 
While  $\theta$  has not coveredged {
   $\theta := \theta - H^{-1}(\theta)\nabla J(\theta)$ 
}
```

Performance of Newton's method:

- ▶ Needs fewer iterations than batch gradient descent

Newton's Method for Optimization

```
Initialize  $\theta$ 
While  $\theta$  has not coveredged {
   $\theta := \theta - \underbrace{H^{-1}(\theta)} \nabla J(\theta)$ 
}
```

Performance of Newton's method:

- ▶ Needs fewer iterations than batch gradient descent
- ▶ Computing $\underbrace{H^{-1}}$ is time consuming

Newton's Method for Optimization

```
Initialize  $\theta$ 
While  $\theta$  has not coveredged {
   $\theta := \theta - H^{-1}(\theta)\nabla J(\theta)$ 
}
```

Performance of Newton's method:

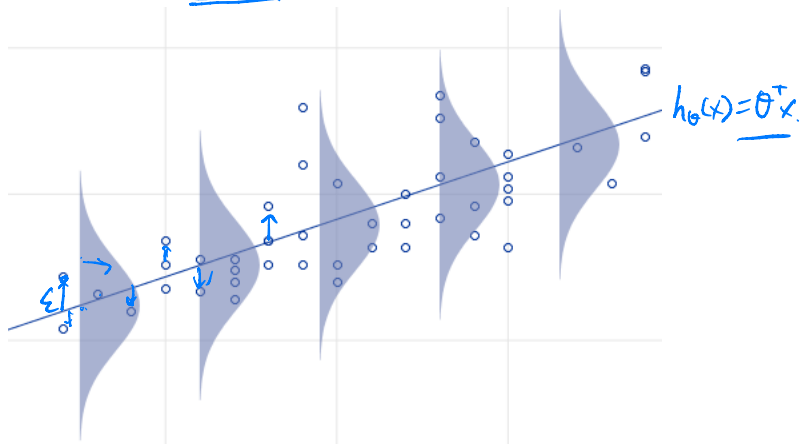
- ▶ Needs fewer iterations than batch gradient descent
- ▶ Computing H^{-1} is time consuming
- ▶ Faster in practice when n is small

Maximum Likelihood Estimation

Consider target y is modeled as

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

and $\epsilon^{(i)}$ are *independently and identically distributed (IID)* to Gaussian distribution $\mathcal{N}(0, \sigma^2)$



Maximum Likelihood Estimation

Consider target y is modeled as

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

and $\epsilon^{(i)}$ are *independently and identically distributed (IID)* to Gaussian distribution $\mathcal{N}(0, \sigma^2)$, then

$$\underline{p(\epsilon^{(i)})} =$$

Maximum Likelihood Estimation

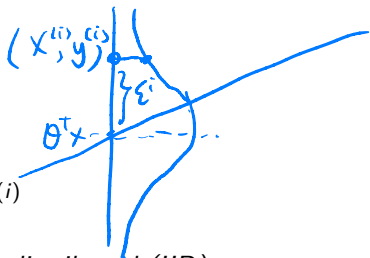
Consider target y is modeled as

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

and $\epsilon^{(i)}$ are *independently and identically distributed (IID)* to Gaussian distribution $\mathcal{N}(0, \sigma^2)$, then

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

$$e^{-\frac{\epsilon^{(i)2}}{2\sigma^2}}$$



Maximum Likelihood Estimation

Consider target y is modeled as

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

and $\epsilon^{(i)}$ are independently and identically distributed (IID) to Gaussian distribution $\mathcal{N}(0, \sigma^2)$, then

$$\underbrace{x^{(i)}, y^{(i)}} \quad p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\epsilon^{(i)2}}{2\sigma^2}\right) \quad \begin{matrix} \epsilon^{(i)} = y^{(i)} - \theta^T x^{(i)} \\ \downarrow \end{matrix}$$
$$\underbrace{p(y^{(i)}|x^{(i)}; \theta)} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

Maximum Likelihood Estimation

The **likelihood** of this model with respect to θ is

$$L(\theta) = \underbrace{p(\vec{y}|X; \theta)}_{\substack{\text{label} \\ \downarrow}} = \prod_{i=1}^m \underbrace{p(y^{(i)}|x^{(i)}; \theta)}_{\substack{\text{input} \\ \downarrow}} \quad \text{due to i.i.d.}$$

Maximum Likelihood Estimation

The **likelihood** of this model with respect to θ is

$$L(\theta) = p(\vec{y}|X; \theta) = \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \theta)$$

Maximum likelihood estimation of θ :

$$\underline{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} L(\theta)$$

Maximum Likelihood Estimation

We compute log likelihood,

$$\begin{aligned} \max_{\theta} \log L(\theta) &= \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi b^2}} + \left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2b^2}\right) \\ &= \underbrace{m \log \frac{1}{\sqrt{2\pi b^2}}}_{\text{constant}} - \frac{1}{b^2} \underbrace{\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2}_{\min_{\theta} J(\theta)} \end{aligned}$$

Maximum Likelihood Estimation

We compute log likelihood,

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

Maximum Likelihood Estimation

We compute log likelihood,

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

Maximum Likelihood Estimation

We compute log likelihood,

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

$$= m \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

Then $\operatorname{argmax}_{\theta} \log L(\theta)$ \equiv $\operatorname{argmin}_{\theta} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$.

Maximum Likelihood Estimation

We compute log likelihood,

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

$$= m \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

Then $\operatorname{argmax}_{\theta} \log L(\theta) \equiv \operatorname{argmin}_{\theta} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$.

i.i.d., $\epsilon^i \sim N(0, \sigma^2)$

Under the assumptions on $\epsilon^{(i)}$, least-squares regression corresponds to the maximum likelihood estimate of θ .

gradient ascent

$$\theta = \theta + \alpha \nabla \log L(\theta)$$

log-likelihood

Linear Regression Summary

How to estimate model parameters θ (or w and b) from data?

- ▶ Least square regression (geometry approach)
- ▶ Maximum likelihood estimation (probabilistic modeling approach)

Linear Regression Summary

How to estimate model parameters θ (or w and b) from data?

- ▶ Least square regression (geometry approach)
- ▶ Maximum likelihood estimation (probabilistic modeling approach)

Other estimation methods exist, e.g. Bayesian estimation

Linear Regression Summary

How to estimate model parameters θ (or w and b) from data?

- ▶ Least square regression (geometry approach)
- ▶ Maximum likelihood estimation (probabilistic modeling approach)

Other estimation methods exist, e.g. Bayesian estimation

How to solve for solutions ?

- ▶ normal equation (close-form solution)
- ▶ gradient descent
- ▶ newton's method

} iterative

Logistic Regression

A binary classification problem

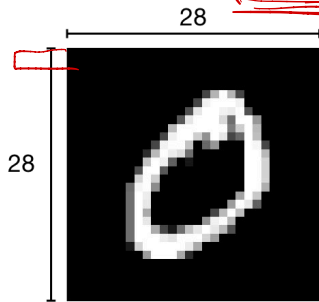
Classify binary digits

- ▶ Training data: 12600 grayscale images of handwritten digits

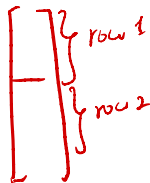


- ▶ Each image is represent by a vector $x^{(i)}$ of dimension $28 \times 28 = \underline{784}$
- ▶ Vectors $x^{(i)}$ are normalized to $[0,1]$

$(\underline{255}, \underline{255}, \underline{255}) \rightarrow$
 $(\underline{0-255})$



$\underline{28 \times 28}$
 $\underline{784}$



A binary classification problem

Classify binary digits

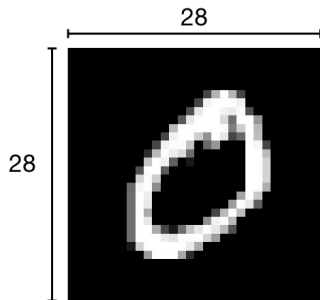
- ▶ Training data: 12600 grayscale images of handwritten digits



- ▶ Each image is represented by a vector $x^{(i)}$ of dimension $28 \times 28 = 784$
- ▶ Vectors $x^{(i)}$ are normalized to $[0,1]$

Binary classification: $\mathcal{Y} = \{0, 1\}$

- ▶ negative class: $y^{(i)} = 0$
- ▶ positive class: $y^{(i)} = 1$

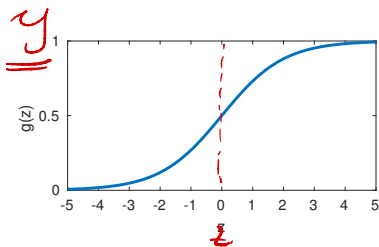


Logistic Regression Hypothesis Function

Sigmoid function

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

- ▶ $g : \mathbb{R} \rightarrow (0, 1)$
- ▶ $\underline{g'(z)} = \underline{g(z)(1 - g(z))}$

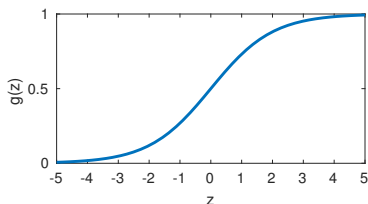


Logistic Regression Hypothesis Function

Sigmoid function

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

- ▶ $g : \mathbb{R} \rightarrow (0, 1)$
- ▶ $g'(z) = g(z)(1 - g(z))$

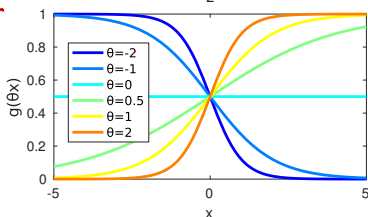
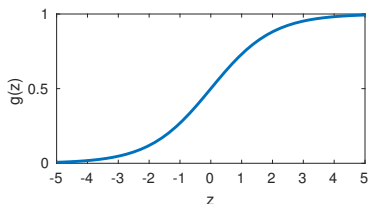


Logistic Regression Hypothesis Function

Sigmoid function

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

- ▶ $g : \mathbb{R} \rightarrow (0, 1)$
- ▶ $g'(z) = g(z)(1 - g(z))$



Hypothesis function for logistic regression:

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

Review: Bernoulli Distribution

A discrete probability distribution of a binary random variable
 $x \in \{0, 1\}$:

$$\begin{aligned} \underline{p(x)} &= \left. \begin{cases} \lambda & \text{if } x = 1 \\ 1 - \lambda & \text{if } x = 0 \end{cases} \right\} \\ &= \underline{p^x (1-p)^{1-x}} = \lambda^x (1-\lambda)^{1-x} \end{aligned}$$

biased.

$$P(\underline{\text{head}}) = \lambda$$

$$\underline{P(\underline{\text{tail}}) = 1 - \lambda}$$



Maximum likelihood estimation for logistic regression

Logistic regression assumes $y|x$ is **Bernoulli distributed**. λ

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

Maximum likelihood estimation for logistic regression

Logistic regression assumes $y|x$ is **Bernoulli distributed**.

- ▶ $p(y = 1 | x; \theta) = h_{\theta}(x)$
- ▶ $p(y = 0 | x; \theta) = 1 - h_{\theta}(x)$

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

Maximum likelihood estimation for logistic regression

Logistic regression assumes $y|x$ is **Bernoulli distributed**.

- ▶ $p(y = 1 | x; \theta) = h_{\theta}(x)$
- ▶ $p(y = 0 | x; \theta) = 1 - h_{\theta}(x)$

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

Given m **independently generated** training examples, the likelihood function is:

$$L(\theta) = p(\vec{y} | X; \theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

Handwritten notes: A red arrow points from the word "log" above the product to the exponent y in the term $h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$. Another red arrow points from the term $h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$ to the right.

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

Handwritten notes: Red underlines are present under $\log(L(\theta))$, $\log h_{\theta}(x^{(i)})$, and $\log(1 - h_{\theta}(x^{(i)}))$.

Maximum likelihood estimation for logistic regression


Logistic regression assumes $y|x$ is **Bernoulli distributed**.

▶ $p(y = 1 | x; \theta) = h_{\theta}(x)$

▶ $p(y = 0 | x; \theta) = 1 - h_{\theta}(x)$

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

Given m **independently generated** training examples, the likelihood function is:


$$L(\theta) = p(\vec{y}|X; \theta) = \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \theta)$$

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

$l(\theta)$ is concave!

Maximum likelihood estimation for logistic regression

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

Solve $\operatorname{argmax}_{\theta} l(\theta)$ using gradient ascent:

$$\frac{\nabla_{\theta} l(\theta)}{=} = \begin{bmatrix} \frac{\partial}{\partial \theta_0} l(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_j} l(\theta) \end{bmatrix}$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} \frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)}) + \frac{1 - y^{(i)}}{1 - h_{\theta}(x^{(i)})} \frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)})$$

$$= \sum_{i=1}^m \left[y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h_{\theta}(x^{(i)})} \right] \frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)})$$

$$= \sum_{i=1}^m \left[y^{(i)} (1 - h_{\theta}(x^{(i)})) x_j^{(i)} - (1 - y^{(i)}) h_{\theta}(x^{(i)}) x_j^{(i)} \right]$$

$$\frac{\partial}{\partial \theta_j} g(\theta^T x^{(i)})$$

$$= g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}$$

$$= h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$= \sum_{i=1}^m \left[y^{(i)} (1 - h_{\theta}(x^{(i)})) - (1 - y^{(i)}) h_{\theta}(x^{(i)}) \right] x_j^{(i)}$$

$$= \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Maximum likelihood estimation for logistic regression

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

Solve $\operatorname{argmax}_{\theta} l(\theta)$ using gradient ascent:

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

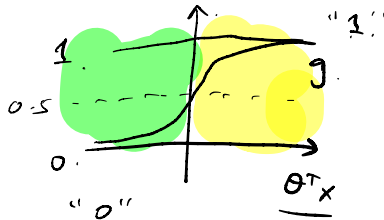
Stochastic Gradient Ascent

```
Repeat until convergence {
  for  $i = 1 \dots m$  {
     $\theta_j = \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$  for every  $j = 1, \dots, n$ .
  }
}
```

$\nabla l(\theta)$

- ▶ Update rule has the same form as least square regression, but with different hypothesis function h_{θ}

Binary Digit Classification



Using the learned classifier

Given an image x , the predicted label is

$$\hat{y} = \begin{cases} 1 & g(\underline{\theta}^T x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Binary digit classification results

	sample size	accuracy
Training	<u>16200</u>	100%
Testing	<u>1225</u>	100%

- ▶ Testing accuracy is 100% since this problem is relatively easy.

Multi-Class Classification

Multiple Binary Classifiers

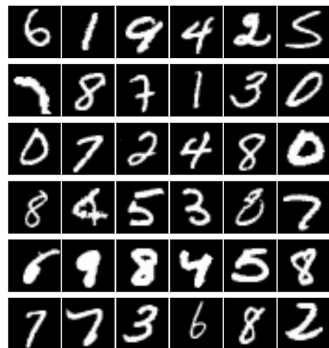
Softmax Regression

Multi-class classification

Each data sample belong to one of $k > 2$ different classes.

$$\mathcal{Y} = \{1, \dots, k\}$$

MNIST Samples



Given new sample $x \in \mathbb{R}^k$, predict which class it belongs.

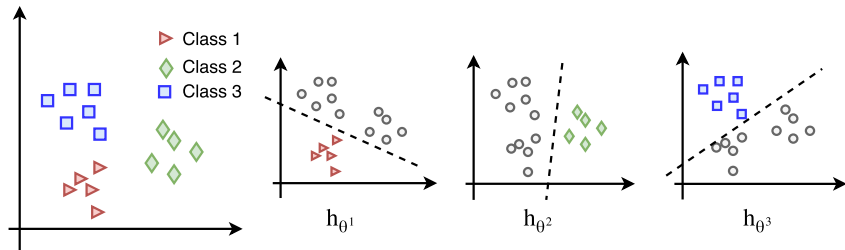
Naive Approach: Convert to binary classification

One-Vs-Rest

Learn k classifiers h_1, \dots, h_k . Each h_i classify one class against the rest of the classes.

Given a new data sample x , its predicted label \hat{y} :

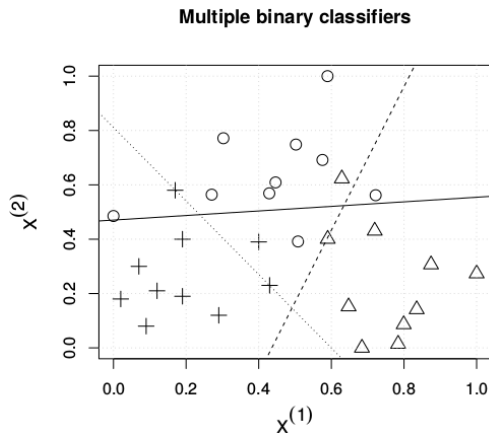
$$\hat{y} = \underset{i}{\operatorname{argmax}} h_i(x)$$



Multiple binary classifiers

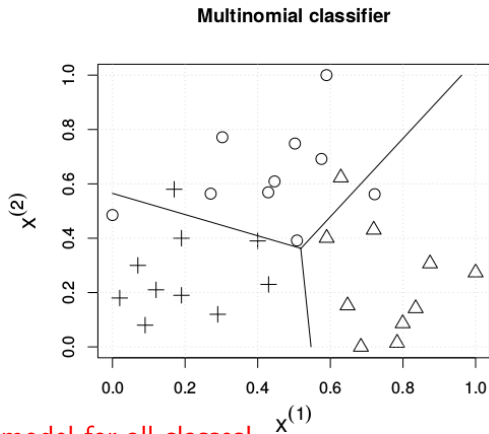
Drawbacks of One-Vs-Rest:

- ▶ Class unbalance: more negative samples than positive samples
- ▶ Different classifiers may have different confidence scales



Drawbacks of One-Vs-Rest:

- ▶ Class imbalance: more negative samples than positive samples
- ▶ Different classifiers may have different confidence scales



Learn one model for all classes!

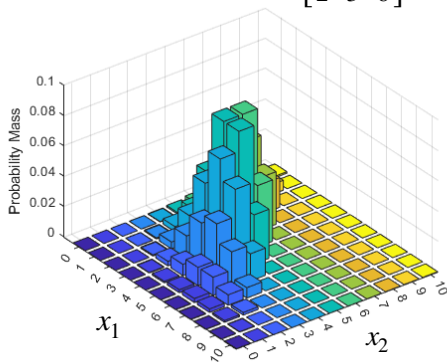
Review: Multinomial Distribution

Models the probability of counts for each side of a k -sided die rolled m times, each side with independent probability ϕ_i



$$\phi_1 + \cdots + \phi_k = 1$$

$$k = 3, n = 10 \quad \phi = \left[\frac{1}{2}, \frac{1}{3}, \frac{1}{6} \right]$$



Extend logistic regression: Softmax Regression

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Extend logistic regression: Softmax Regression

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Hypothesis function for sample x :

$$h_{\theta}(x) = \begin{bmatrix} p(y = 1|x; \theta) \\ \vdots \\ p(y = k|x; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x_j}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_k^T x} \end{bmatrix} = \text{softmax}(\theta^T x)$$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Extend logistic regression: Softmax Regression

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Hypothesis function for sample x :

$$h_{\theta}(x) = \begin{bmatrix} p(y = 1|x; \theta) \\ \vdots \\ p(y = k|x; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x_j}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_k^T x} \end{bmatrix} = \text{softmax}(\theta^T x)$$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Parameters: $\theta = \begin{bmatrix} - & \theta_1^T & - \\ & \vdots & \\ - & \theta_k^T & - \end{bmatrix}$

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)}) \mathbf{1}_{\{y^{(i)}=l\}}\end{aligned}$$

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)}) \mathbf{1}_{\{y^{(i)}=l\}} \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}_{\{y^{(i)} = l\}} \log p(y^{(i)} = l | x^{(i)})\end{aligned}$$

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)}) \mathbf{1}\{y^{(i)}=l\} \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}\{y^{(i)} = l\} \log p(y^{(i)} = l | x^{(i)}) \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}\{y^{(i)} = l\} \log \frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}}\end{aligned}$$

Softmax Regression

Derive the stochastic gradient descent update:

- ▶ Find $\nabla_{\theta_l} \ell(\theta)$

$$\nabla_{\theta_l} \ell(\theta) = \sum_{i=1}^m \left[\left(\mathbf{1}\{y^{(i)} = l\} - P(y^{(i)} = l | x^{(i)}; \theta) \right) x^{(i)} \right]$$

Property of Softmax Regression

- ▶ Parameters $\theta_1, \dots, \theta_k$ are not independent:
$$\sum_j p(y = j|x) = \sum_j \phi_j = 1$$
- ▶ Knowing $k - 1$ parameters completely determines model.

Invariant to scalar addition

$$p(y|x; \theta) = p(y|x; \theta - \psi)$$

Proof.

Relationship with Logistic Regression

When $K = 2$,

$$h_{\theta}(x) = \frac{1}{e^{\theta_1^T x} + e^{\theta_2^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ e^{\theta_2^T x} \end{bmatrix}$$

Relationship with Logistic Regression

When $K = 2$,

$$h_{\theta}(x) = \frac{1}{e^{\theta_1^T x} + e^{\theta_2^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ e^{\theta_2^T x} \end{bmatrix}$$

Replace $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$ with $\theta_* = \theta - \begin{bmatrix} \theta_2 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \theta_1 - \theta_2 \\ 0 \end{bmatrix}$,

$$\begin{aligned} h_{\theta}(x) &= \frac{1}{e^{\theta_1^T x - \theta_2^T x} + e^{0^T x}} \begin{bmatrix} e^{(\theta_1 - \theta_2)^T x} \\ e^{0^T x} \end{bmatrix} \\ &= \begin{bmatrix} \frac{e^{(\theta_1 - \theta_2)^T x}}{1 + e^{(\theta_1 - \theta_2)^T x}} \\ \frac{1}{1 + e^{(\theta_1 - \theta_2)^T x}} \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{1 + e^{-(\theta_1 - \theta_2)^T x}} \\ 1 - \frac{1}{1 + e^{-(\theta_1 - \theta_2)^T x}} \end{bmatrix} = \begin{bmatrix} g(\theta_*^T x) \\ 1 - g(\theta_*^T x) \end{bmatrix} \end{aligned}$$

When to use Softmax?

- ▶ When classes are mutually exclusive: use Softmax
- ▶ Not mutually exclusive: multiple binary classifiers may be better