

Learning From Data

Lecture 2: Linear Regression & Logistic Regression

Yang Li yangli@sz.tsinghua.edu.cn

September 23, 2022

Outline

Introduction

Today's Lecture

Supervised Learning (Part I)

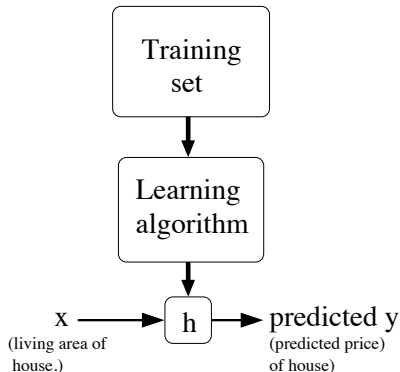
- ▶ Linear Regression
- ▶ Binary Classification
- ▶ Multi-Class Classification

Review: Supervised Learning

- ▶ Input space: \mathcal{X} , Target space: \mathcal{Y}

Review: Supervised Learning

- ▶ Input space: \mathcal{X} , Target space: \mathcal{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

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

Outline

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n = \underline{\underline{\theta^T x}}$$

Linear Regression

Linear Regression Model

Ordinary Least Square

- geometric

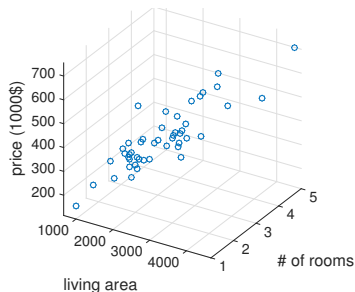
Maximum Likelihood Estimation

- probabilistic

Linear Regression

Example: predict Portland housing price

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



Linear Approximation

A linear model

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

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,

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

Alternative Notation

$$h(x) = w_1x_1 + w_2x_2 + b$$

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

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

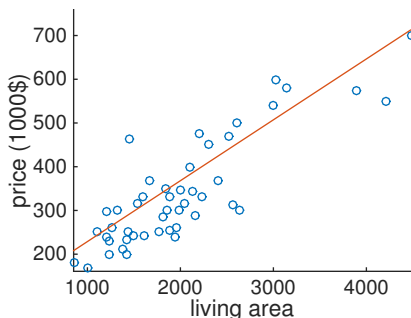
Apply model to new data

Suppose we have the optimal parameters θ , e.g.

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

make a prediction of new feature x :

$$\hat{y} = h_{\theta}(x) = \theta^T x$$



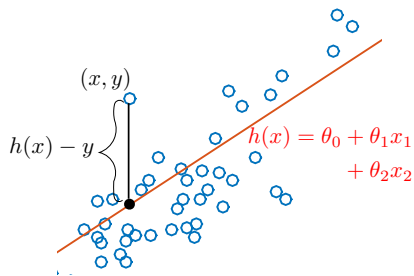
Model Estimation

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

Model Estimation

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

Least Square Estimation

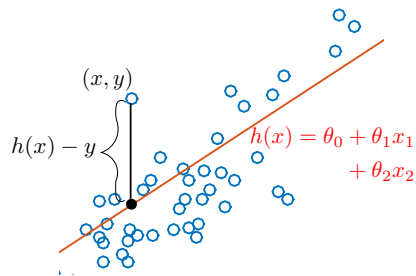


geometric approach

Model Estimation

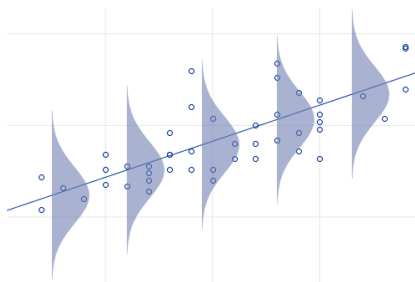
How to estimate model parameters θ (or w and b) from data?

Least Square Estimation



geometric approach

Maximum Likelihood Estimation

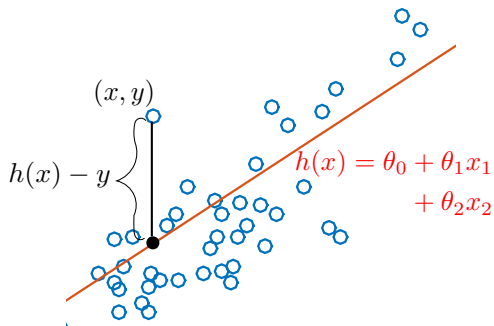


Probabilistic approach

Ordinary Least Square

Cost function:

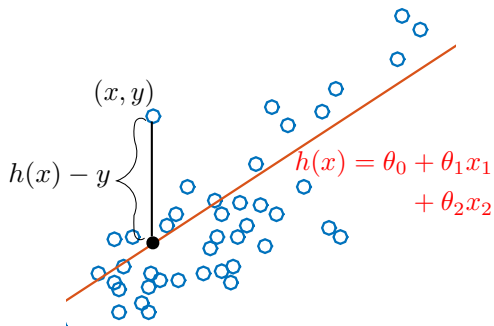
$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$



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:

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

The **ordinary Least square problem** is:

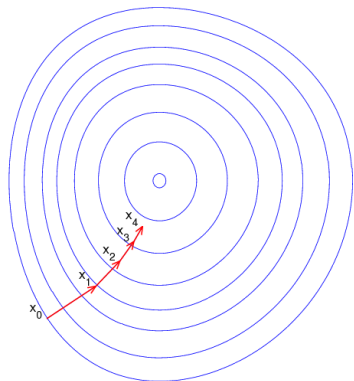
$$\begin{aligned} & \min_{\theta} J(\theta) \\ &= \min_{\theta} \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \end{aligned}$$

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 $J(\theta)$:

$$\theta := \theta - \alpha \nabla J(\theta)$$

α is the **learning rate**

Review: Convex function

Definition (Convex set)

Let S be a vector space, any subset $C \subseteq S$ is **convex** if for any $x, y \in C$, $0 \leq \lambda \leq 1$, affine combination¹ $\lambda x + (1 - \lambda)y \in C$

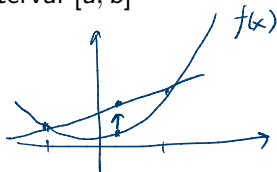
¹An affine combination is a linear combination where coefficients sum to 1.

Definition (Convex function)

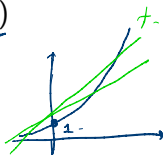
A function $f(x)$ is **convex** on a convex set C if for any $x_1, x_2 \in C$ and $0 \leq \lambda \leq 1$,

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

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



- x^2
- e^x
- $\|x\|_p$
- composition rules.



Definition (Convex function)

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.

$$\theta := \theta - \alpha \nabla J(\theta)$$

For the ordinary least square problem,

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

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

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

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

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

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

$$= \frac{\partial}{\partial \theta_j} \sum_{j=1}^n \theta_j \cdot x_j^{(i)} = \frac{\partial}{\partial \theta_j} \theta_j x_j^{(i)}$$

$$= x_j^{(i)}$$

For the ordinary least square problem,

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

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

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  
}
```

while θ is not converged:

for j in range(n):

$\theta_j = \theta_j + \dots$]

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  
}
```

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 }  
}
```

Stochastic gradient descent (SGD)

```
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 }  
  }  
}
```

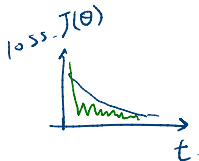
θ is updated each time a training example is read

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  
  }  
}
```



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

$$x^{(i)} \in \mathbb{R}^n \quad (i=1, \dots, m)$$

The matrix notation

$$\underline{X} = \begin{matrix} \underbrace{\hspace{10em}}_n \\ \left[\begin{array}{c} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{array} \right] \end{matrix} \left. \vphantom{\begin{matrix} \underbrace{\hspace{10em}}_n \\ \left[\begin{array}{c} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{array} \right] \end{matrix}} \right\} m \end{matrix}, \quad \underline{\vec{y}} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

X is called the **design matrix**.

Minimize $J(\theta)$ Analytically

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

The matrix notation

$$\underbrace{(X\theta - y)^T (X\theta - y)}_{\|z\|^2} = \sum_{i=1}^m z_i^2 = \sum_{i=1}^n \underbrace{(x^{(i)T} \theta - y^{(i)})^2}$$

$$\underbrace{X}_{m \times n} = \begin{bmatrix} - & (x^{(1)})^T & - \\ - & (x^{(2)})^T & - \\ & \vdots & \\ - & (x^{(m)})^T & - \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

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

Compute the gradient of $J(\theta)$: *least sq. function*

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

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

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

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

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

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

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

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

$$\begin{aligned}X^T X \theta &= X^T y \\ \theta &= \underline{(X^T X)^{-1} X^T y}\end{aligned}$$

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

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

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

The Normal equation

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

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

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

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

Which method to use?

$$\min_{\theta} J(\theta)$$

→ more iteration

gradient descent

iterative solution

need to
tune α ! }
—

- distributed gradient descent

- SGD

normal equation

exact solution

faster small-scale problem
optimal solution

$(X^T X)^{-1}$ → numerical issue

Which method to use?

gradient descent	normal equation
iterative solution	exact solution
need to choose proper learning parameter α for cost function to converge	

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 <u>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 <u>samples m</u>	

Which method to use?

$$\frac{1}{2} \underbrace{(X\theta - y)^T (X\theta - y)}$$

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 m is large

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

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

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

Solving $f(x) = 0$.

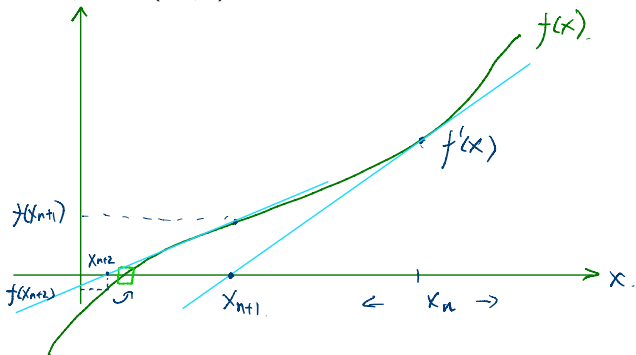
Geometric intuition of Newton's method

At step $n + 1$: $\underline{x_n}, y_n = f(x_n)$.

$$f(x) = f'(x)(x - x_n) + \overset{y_n}{f(x_n)} = 0$$

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

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



Newton's Method Demo

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

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

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

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

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

- ▶ θ is one-dimensional:

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

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

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

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

- ▶ θ is one-dimensional:

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

- ▶ θ is multidimensional:

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

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

a.k.a Newton-Raphson method

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

($n \times n$.)

Newton's Method for Optimization

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


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 {
   $\theta := \theta - \underbrace{H^{-1}(\theta)}_{n \times n} \nabla J(\theta)$ 
}
```

Performance of Newton's method:

- ▶ Needs fewer iterations than batch gradient descent
- ▶ Computing 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

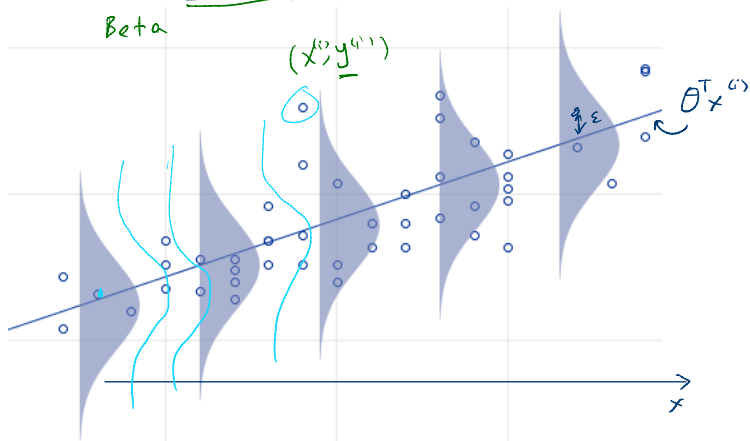
Maximum Likelihood Estimation

Consider target y is modeled as

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

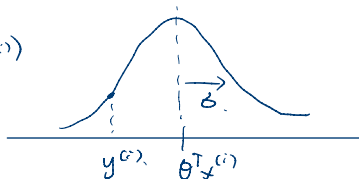
$$E\left[\frac{\hat{\theta}^T x^{(i)} + \epsilon^i}{\epsilon^i \sim \mathcal{N}(0, \sigma^2)}\right] = \theta^T x^i$$

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



Maximum Likelihood Estimation

pdf($y^{(i)}$)



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

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\epsilon^{(i)2}}{2\sigma^2}}$$

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

Given $x^{(i)}, y^{(i)}$.

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

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

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

Maximum Likelihood Estimation

Consider target y is modeled as

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

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

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

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

Maximum Likelihood Estimation

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \quad \underline{p(y_1, \dots, y_m | x; \theta)} = \prod_{i=1}^m \underline{p(y^{(i)} | x^{(i)}; \theta)}$$

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

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

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 :

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

Maximum Likelihood Estimation

We compute log likelihood,

$$\begin{aligned} \max_{\theta} \log L(\theta) &= \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \right) \\ &= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi b^2}} + \left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2b^2} \right) \\ &= m \log \frac{1}{\sqrt{2\pi b^2}} - \frac{1}{2b^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \\ &\quad \underbrace{\hspace{10em}}_{\frac{1}{2} \cdot J(\theta)} \end{aligned}$$

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,

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

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

$\mathcal{J}(\theta)$

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,

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

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

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

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

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

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

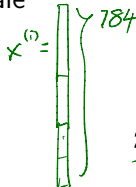
Outline

Logistic Regression

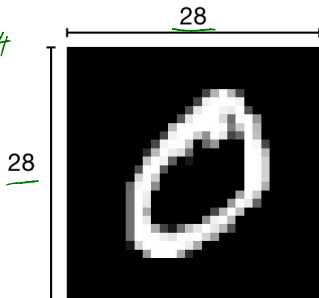
A binary classification problem

Classify binary digits

- ▶ Training data: 12600 grayscale images of handwritten digits



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



A binary classification problem

Classify binary digits

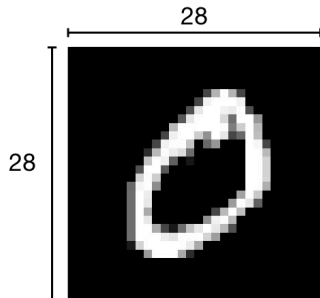
- ▶ Training data: 12600 grayscale images of handwritten digits



- ▶ 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

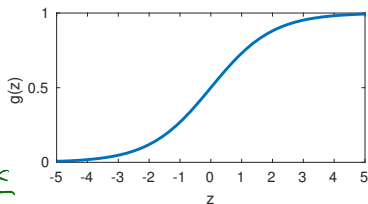
Sigmoid function

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

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

▶ $\underline{g'(z)} = \underline{g(z)}(1 - \underline{g(z)})$

Homework

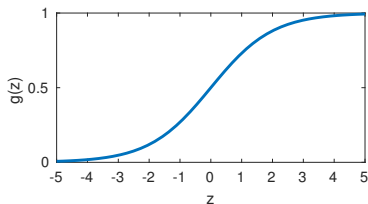


Logistic Regression Hypothesis Function

Sigmoid function

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

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

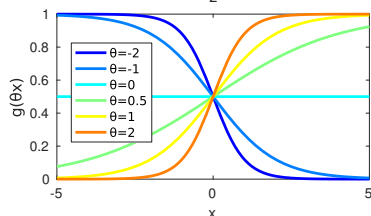
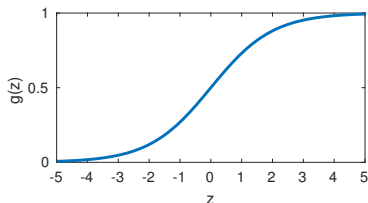


Logistic Regression Hypothesis Function

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:

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

Review: Bernoulli Distribution

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

$$\underline{y^{(i)} | x^{(i)}} \sim \underline{\mathcal{N}(\theta^T x^{(i)}, b^2)}$$

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



Maximum likelihood estimation for logistic regression

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

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

Maximum likelihood estimation for logistic regression

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

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

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

$$p(x; \lambda) = \lambda^x (1 - \lambda)^{1-x}$$

↓

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

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:

$$y \log h_{\theta}(x) + (1-y) \log(1 - h_{\theta}(x))$$

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

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

Maximum likelihood estimation for logistic regression

① f is concave iff $-f$ is convex ②



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:

$$\begin{aligned} h_{\theta}(x^{(i)}) &= g(\theta^T x^{(i)}) \\ \frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)}) &= g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)} \\ &= h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x_j^{(i)} \end{aligned}$$
$$\begin{aligned} \frac{\partial l(\theta)}{\partial \theta_j} &= \sum_{i=1}^m y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_j} + \frac{1 - y^{(i)}}{1 - h_{\theta}(x^{(i)})} (-1) \cdot \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_j} \\ &= \sum_{i=1}^m [y^{(i)} (1 - h_{\theta}(x^{(i)})) x_j^{(i)} - (1 - y^{(i)}) h_{\theta}(x^{(i)}) x_j^{(i)}] \\ &= \sum_{i=1}^m [y^{(i)} (1 - h_{\theta}(x^{(i)})) - (1 - y^{(i)}) h_{\theta}(x^{(i)})] x_j^{(i)} \\ &= \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \\ &\quad g(\theta^T x^{(i)}) \end{aligned}$$

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  $j$   
  }  
}
```

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

Binary Digit Classification

$$\underline{\theta^T x}$$

Using the learned classifier

Given an image x , the predicted label is

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

Binary digit classification results

	sample size	accuracy
Training	16200	100%
Testing	1225	100%

► Testing accuracy is 100% since this problem is relatively easy.

Q: label imbalance ?

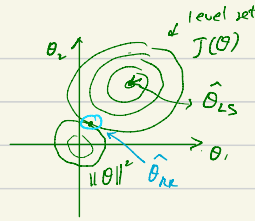
focal loss

Homework Preview

① Ridge regression

$$\min_{\theta} J(\theta)$$

$$J(\theta) = \underbrace{\|y - X\theta\|^2}_{+ (\theta^T \theta)} + \lambda \|\theta\|^2$$



- underdetermined $X\theta = y$ $m < n$
- prevent overfitting
- highly correlated variables. $(X^T X)^{-1} \leftarrow$ numerical issues

$$\hat{\theta}_{RR} = (X^T X + \lambda I)^{-1} X^T y$$

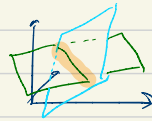
$$\text{When } \lambda \rightarrow 0, \hat{\theta}_{RR} \rightarrow \hat{\theta}_{LS}$$

② over-parameterized model:

Implicit Bias of Gradient Descent:

Given an under-determined system

$$y = \underbrace{Ax}_{\text{full rank}} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$



assuming A has full row rank.

$$\min_x \|y - Ax\|$$

GD: Let $x_0 = 0$

$$\text{update: } x_k = x_{k-1} - \alpha \nabla f(x)$$

$$x^\infty \rightarrow x_{LN}$$

Definition Least Norm Solution of a Linear System:

$$x_{LN} = \arg \min_x \|x\|$$

$$\text{st. } Ax = y$$

$$x_{LN} = A^T (AA^T)^{-1} y$$

↑ not the same as $(A^T A)^{-1} A^T y$

Outline

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 to.

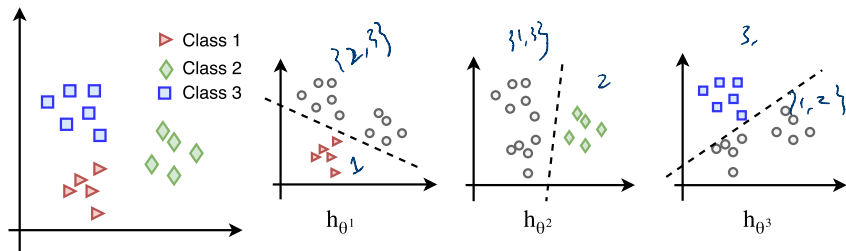
Naive Approach: Convert to binary classification

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} :

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



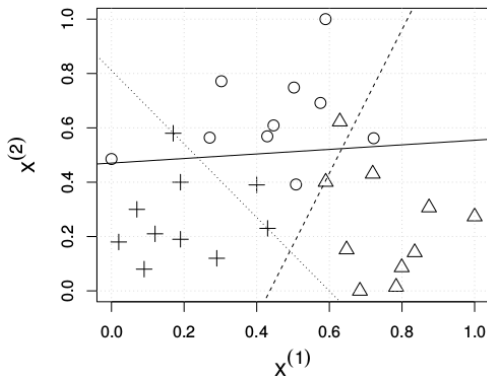
Multiple binary classifiers

Drawbacks of One-Vs-Rest:

*100 samples, 20 classes
5 positive samples for each h_i*

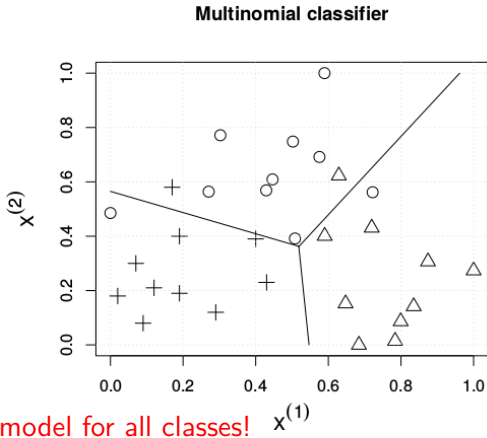
- ▶ Class unbalance: more negative samples than positive samples
- ▶ Different classifiers may have different confidence scales

Multiple binary classifiers



Drawbacks of One-Vs-Rest:

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



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

3-sided die.



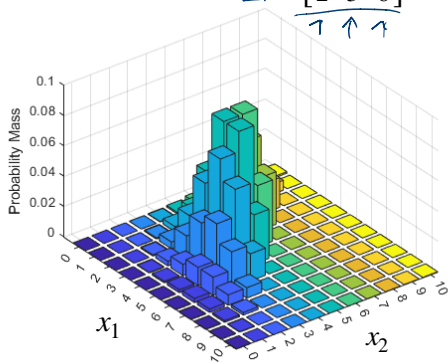
$$\phi_1 + \dots + \phi_k = 1$$

$$k = 3, n = 10$$

$$\underline{\phi} = \left[\frac{1}{2}, \frac{1}{3}, \frac{1}{6} \right]$$

↑ ↑ ↑

$$\phi = \left\{ \underbrace{\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6}} \right\}$$



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 :

$$\underbrace{h_\theta(x)}_{\in \mathbb{R}^k} = \begin{bmatrix} p(y=1|x; \theta) \\ p(y=2|x; \theta) \\ \vdots \\ p(y=k|x; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x_j}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_k^T x} \end{bmatrix} = \underline{\text{softmax}(\theta^T x)}$$

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

Extend logistic regression: Softmax Regression

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

Hypothesis function for sample x :

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

$x \in \mathbb{R}^n$

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

Parameters: $\theta = \underbrace{\begin{bmatrix} - & \theta_1^T & - \\ & \vdots & \\ - & \theta_k^T & - \end{bmatrix}}_{\mu}$ } k .

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

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

samples \rightarrow m k \leftarrow *classes*

Softmax Regression

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

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

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