# Learning From Data
# Lecture 2: Linear Regression & Logistic Regression

Yang Li    yangli@sz.tsinghua.edu.cn

March 7, 2024

# Outline

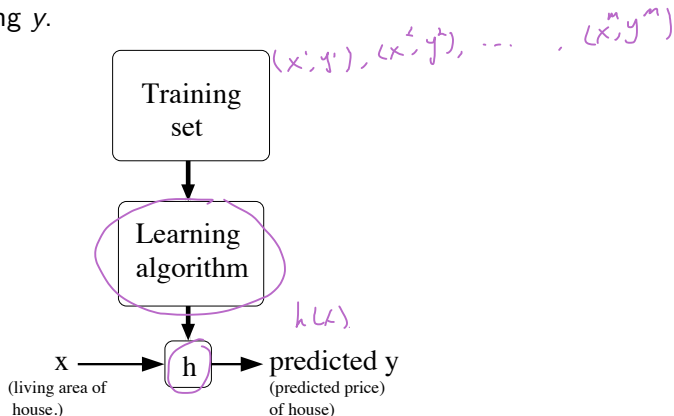Introduction

# 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 $h : \mathcal{X} \to \mathcal{Y}$ so that $h(x)$ is a "good" predictor for the corresponding $y$.

$(x^1, y^1), (x^2, y^2), \ldots , (x^m, y^m)$

```
        ┌─────────────┐
        │  Training   │
        │    set      │
        └─────────────┘
               │
               ▼
        ┌─────────────┐
        │  Learning   │
        │  algorithm  │
        └─────────────┘
               │
               ▼  h(x)
  x ──────▶  ( h ) ──────▶ predicted y
(living area of              (predicted price)
 house.)                     of house)
```

# Review: Supervised Learning

- $y$ is discrete (categorical): **classification problem**
- $y$ is continuous (real value): **regression problem** $\Leftarrow$

# Outline

$h(x)$ is a linear function

## Linear Regression

# 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

*bias.*

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

$\theta_i$'s are called **parameters**.

# Linear Approximation

A linear model

bias

weights

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

$\theta_i$'s are called **parameters**.

Using vector notation,

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

$$= \theta_0 \cdot 1 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2$$

# Alternative Notation

$$h(x) = w_1 x_1 + w_2 x_2 + b$$

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

$$h(x) = \underline{\theta^T x}$$

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

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \qquad b \in \mathbb{R}$$
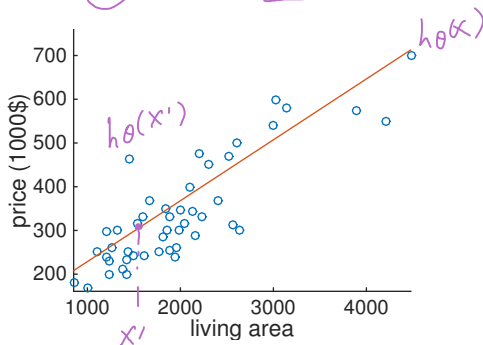
# Apply model to new data

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

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

make a prediction of new feature $x$:
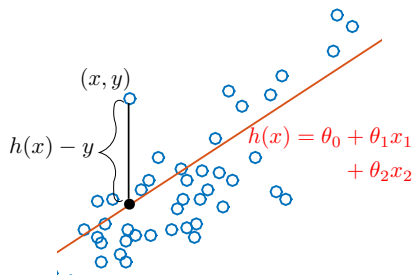
$$\hat{y} = h_\theta(x) = \theta^T x$$

# Model Estimation

How to estimate model parameters $\theta$ (or $w$ and $b$) from data?

# Model Estimation

How to estimate model parameters $\theta$ (or $w$ and $b$) from data?
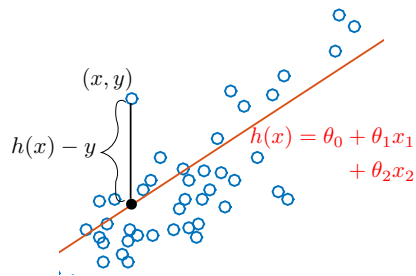
## Least Square Estimation



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

geometric approach

# Model Estimation

How to estimate model parameters $\theta$ (or $w$ and $b$) from data?
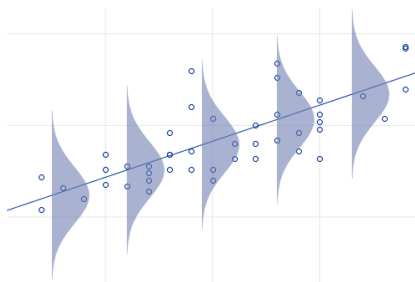
## Least Square Estimation



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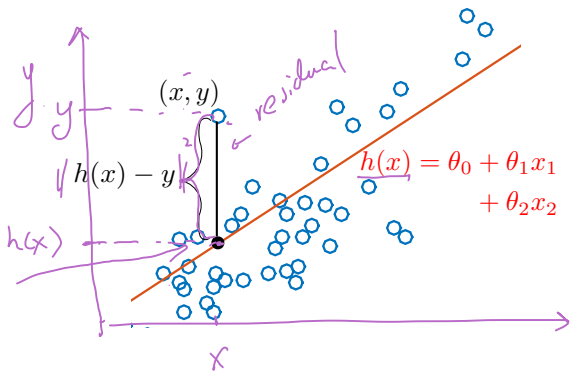
geometric approach

## Maximum Likelihood Estimation



Probabilistic approach

# Ordinary Least Square

Cost function:

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



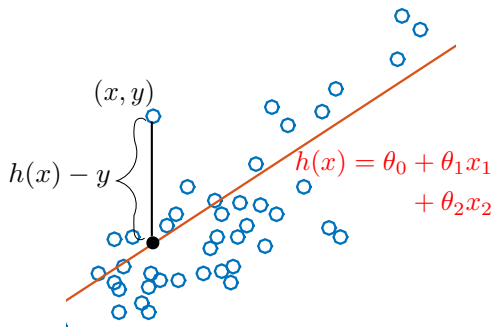$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

# 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:

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

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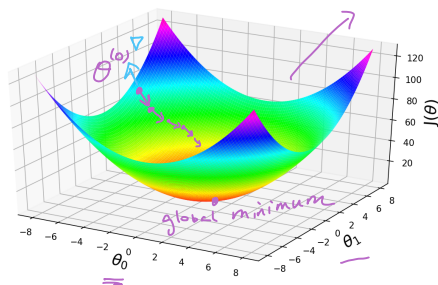
$$\min_{\theta} J(\theta)$$

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

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, $\theta^{(o)}$ repeatedly change $\theta$ to decrease $J(\theta)$:

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

$\alpha$ 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 $x, y \in C$, $0 \le \lambda \le 1$, affine combination[1] $\lambda x + (1 - \lambda)y \in C$

---

[1]An 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 $0 \le \lambda \le 1$,

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

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

### 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.

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

For the ordinary least square problem,

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

$m$ samples

$\theta \in \mathbb{R}^n.$

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

$j$th

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

$\theta^T x^{(i)} = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \cdots \theta_n x_n^{(i)}$

$\cdot$ if $j = 2$, then only $\theta_j x_j^{(i)}$ depends on $\theta_j$.

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

$x_j^{(i)}$.

For the ordinary least square problem,
$J(\theta) = \frac{1}{2} \sum_{i=1}^{m} (h(x^{(i)}) - y^{(i)})^2 = \frac{1}{2} \sum_{i=1}^{m} (\theta^T x^{(i)} - y^{(i)})^2$,

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

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

# Gradient descent for ordinary least square

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

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`
```
}
```

$\text{for} \quad j = 1 \cdots n :$
$\triangle = 0$
$\quad \begin{bmatrix} \text{for} \quad i = 1 \cdots m : \\ \quad \triangle = \triangle + (y^i - h_\theta(x^i)) x_j^i \end{bmatrix}$
$\theta_j := \theta_j + \alpha \cdot \triangle$

# Gradient descent for ordinary least square

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

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

## Batch Gradient Descent

```
Repeat until convergence{
    θⱼ = θⱼ + α ∑ᵢ₌₁ᵐ(y⁽ⁱ⁾ − hθ(x⁽ⁱ⁾))xⱼ⁽ⁱ⁾ for every j
}
```

$\theta$ is only updated after we have seen all $m$ training samples.

## Batch gradient descent

```
Repeat until convergence{
    θⱼ = θⱼ + α ∑ᵢ₌₁ᵐ (y⁽ⁱ⁾ − hθ(x⁽ⁱ⁾))xⱼ⁽ⁱ⁾  for every j
}
```

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

*outer loop.*

## Stochastic gradient descent

```
Repeat until convergence{
    for i = 1...m {
        θⱼ = θⱼ + α(y⁽ⁱ⁾ − hθ(x⁽ⁱ⁾))xⱼ⁽ⁱ⁾  for every j
    }
}
```

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

$\theta$ is updated each time a training example is read

## Batch gradient descent

```
Repeat until convergence{
    θ_j = θ_j + α  ∑_{i=1}^m  (y^{(i)} − h_θ(x^{(i)}))x_j^{(i)}  for every j
}
```

## Stochastic gradient descent

```
Repeat until convergence{
    for i = 1...m {
        θ_j = θ_j + α(y^{(i)} − h_θ(x^{(i)}))x_j^{(i)}  for every j
    }
}
```

$\theta$ is updated each time a training example is read

- ▶ Stochastic gradient descent gets $\theta$ close to minimum much faster (Video)
- ▶ Good for regression on large data

# Minimize $J(\theta)$ Analytically

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

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

The matrix notation

$$x_1^{(1)}\ x_2^{(1)}\ \cdots\ x_n^{(1)} \qquad \theta$$

$$X = m\left\{\begin{bmatrix} -\ (x^{(1)})^T\ - \\ -\ (x^{(2)})^T\ - \\ \vdots \\ -\ (x^{(m)})^T\ - \end{bmatrix}\right. , \begin{bmatrix} \theta_1 \\ \\ \theta_n \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}^{m\times 1}$$

$X$ is called the **design matrix**.

# Minimize $J(\theta)$ Analytically

The matrix notation

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

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

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

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

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

① $\dfrac{\partial x^T A x}{\partial x} = Ax + A^T x$
$\qquad = (A + A^T)x$

if $A$ is symmetric

$\qquad A = A^T$;

$\dfrac{\partial x^T A x}{\partial x} = 2Ax.$

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

$\qquad = \frac{1}{2} (\theta^T X^T X\theta - \underbrace{\theta^T X^T y}_{\in \mathbb{R}} - \underbrace{y^T X\theta}_{\in \mathbb{R}} + y^T y)$

② $\dfrac{\partial x^T a}{\partial x} = \underline{a}$

$\qquad = \frac{1}{2}\theta^T X^T X\theta - \frac{1}{2} 2\, \theta^T X^T y + \frac{1}{2} y^T y.$

$\nabla_\theta J(\theta) = \frac{\partial}{\partial \theta} J(\theta) = \frac{1}{2} \cdot \frac{\partial (\theta^T X^T X \theta)}{\partial \theta} - \frac{\partial}{\partial \theta} \theta^T X^T y + \frac{\partial y^T y}{\partial \theta}$  $0.$

$X^T X \theta - X^T y = 0.$

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

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 $J(\theta)$ :

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

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

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

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

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

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

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

The Normal equation

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

col(X)

if $X^T X$ is not full rank,
solution $\theta^*$ is not unique.

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

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

$$= X^T X \theta - X^T y$$

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

When is $X^T X$ not invertible?

The Normal equation

① features are not independent.

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

② # of linearly independent samples are fewer than # of features.

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

$X = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$ (2×3)

$X^T X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

$\text{rank}(X^T X) = 2$

rank deficiency.

$Ax = y$

$(X^T X)\theta = X^T y$.

# Which method to use?

numerical            analytical

| **gradient descent** | **normal equation** |
| --- | --- |
| iterative solution | exact solution |
| | |
| | |

# Which method to use?

| gradient descent | normal equation |
| --- | --- |
| iterative solution | exact solution |
| need to choose proper learning parameter $\alpha$ for cost function to converge | |

# Which method to use?

| **gradient descent** | **normal equation** |
| --- | --- |
| iterative solution | exact solution |
| need to choose proper learning parameter $\alpha$ for cost function to converge | numerically unstable when $X$ is ill-conditioned. e.g. features are highly correlated |

$X^T X$

# Which method to use?

| gradient descent | normal equation |
|---|---|
| iterative solution | exact solution |
| need to choose proper learning parameter $\alpha$ 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 $\alpha$ 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 $\theta$ in $\nabla_\theta J(\theta) = 0$

## Newton's method

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

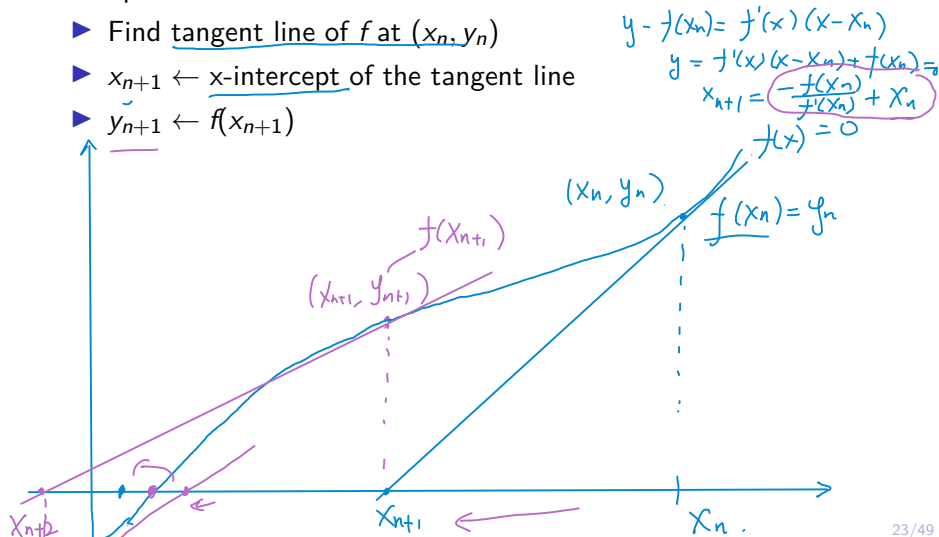▶ Start an initial guess $x^{(0)}$

▶ 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 underline{tangent line of $f$} at $(x_n, y_n)$

▶ $x_{n+1} \leftarrow$ x-intercept of the tangent line

▶ $y_{n+1} \leftarrow f(x_{n+1})$

$y - f(x_n) = f'(x)(x - x_n)$

$y = f'(x)(x - x_n) + f(x_n) = 0$

$x_{n+1} = \boxed{-\dfrac{f(x_n)}{f'(x_n)} + x_n}$

$f(x) = 0$

$(x_n, y_n)$

$f(x_n) = y_n$

$f(x_{n+1})$

$(x_{n+1}, y_{n+1})$

$x_{n+2}$   $x_{n+1}$   $x_n$

# Newton's Method Demo

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

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

$x = x - \dfrac{f(x)}{f'(x)}$

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

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

- $\theta$ is one-dimensional:

$$f = J'(\theta) = \nabla_\theta J(\theta).$$

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

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

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

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

- $\theta$ is one-dimensional:

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

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

$(n \times n)$

- $\theta$ is multidimensional: $\theta \in \mathbb{R}^n$

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

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

a.k.a Newton-Raphson method

# Newton's Method for Optimization

```
Initialize θ
While θ has not coverged {
    θ := θ − H⁻¹(θ)∇J(θ)
}
```

# Newton's Method for Optimization

```
Initialize θ
While θ has not coverged {
   θ := θ − H⁻¹(θ)∇J(θ)
}
```

Performance of Newton's method:

▶ Needs fewer interations than batch gradient descent

# Newton's Method for Optimization

```
Initialize θ
While θ has not coverged {
    θ := θ − H⁻¹(θ)∇J(θ)
}
```

Performance of Newton's method:

▶ Needs fewer interations than batch gradient descent ◡̈

▶ Computing $H^{-1}$ is time consuming ⌣̀

# Newton's Method for Optimization

```
Initialize θ
While θ has not coverged {
    θ := θ − H⁻¹(θ)∇J(θ)
}
```

Performance of Newton's method:
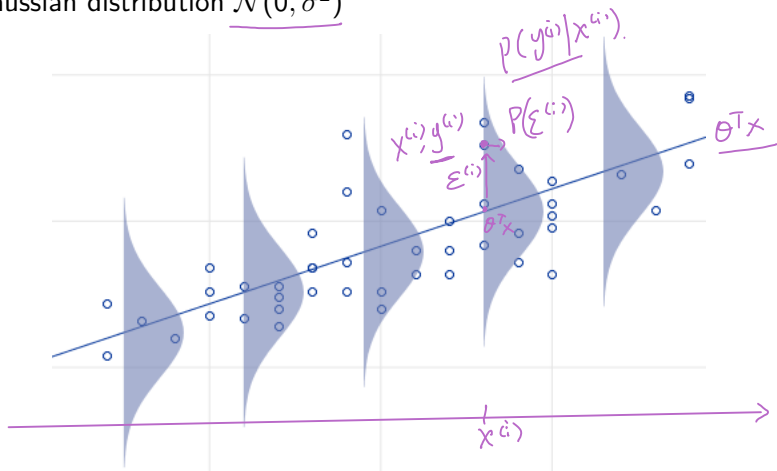
► Needs fewer interations than batch gradient descent

► 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

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

# Maximum Likelihood Estimation

Consider target $y$ is modeled as

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

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

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

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

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

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

*likelihood of one sample*

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

$\varepsilon^{(i)} \,/\quad y^{(i)}|x^{(i)} \sim \underline{i.i.d.}\ \mathcal{N}(0, b^2)$

The **likelihood** of this model with respect to $\theta$ is

$$L(\theta) = \underbrace{p(\vec{y}|X; \theta)}_{\substack{(m) \quad (m \times n)}} = \prod_{i=1}^{m} \underline{p(y^{(i)}|x^{(i)}; \theta)}$$

# Maximum Likelihood Estimation

The **likelihood** of this model with respect to $\theta$ is

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

**Maximum likelihood estimation of** $\theta$:

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}}\, L(\theta)$$

# Maximum Likelihood Estimation

We compute log likelihood,

$$\log L(\theta) = \log \left( \prod_{i=1}^{m} p(y^{(i)}|x^{(i)}; \theta) \right) = \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)}; \theta)$$

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

# Maximum Likelihood Estimation

We compute log likelihood,

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

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

# Maximum Likelihood Estimation

We compute log likelihood,

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

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

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

# Maximum Likelihood Estimation

We compute log likelihood,

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

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

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

$J(\theta)$ in ordinary least square

doesn't depend on $\theta$

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,

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

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

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

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

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

# Linear Regression Summary

How to estimate model parameters $\theta$ (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 $\theta$ (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 $\theta$ (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*  MAP.

How to solve for solutions ?

- ▶ normal equation (close-form solution)  analytical

- ▶ gradient descent  ⎫
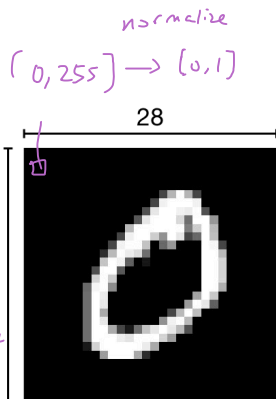- ▶ newton's method  ⎬ numerical

Logistic Regression
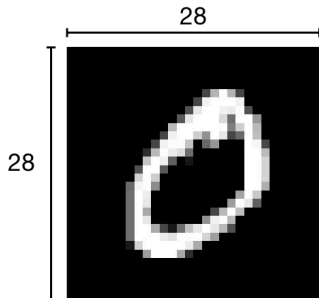
# A binary classification problem

## 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 $[0,1]$

normalize

$[0, 255] \longrightarrow [0, 1]$

28

28



28

28

784

784

$x^{(i)}$.

# 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]

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

- negative class: $y^{(i)} = 0$
- positive class: $y^{(i)} = 1$
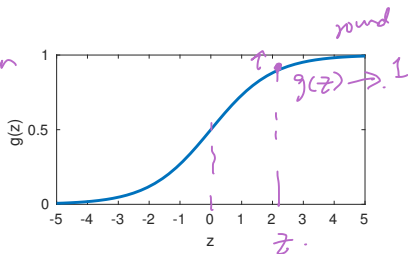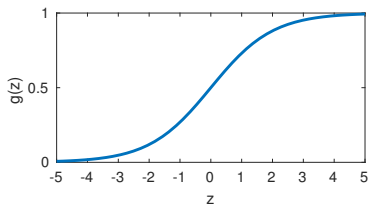
# Logistic Regression Hypothesis Function

### Sigmoid function

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

→ parameter

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

$$\frac{\partial}{\partial z}\left(\frac{1}{1 + e^{-z}}\right)$$



round

$g(z) \to 1$

$z$.

# Logistic Regression Hypothesis Function

Sigmoid function

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

- $g : \mathbb{R} \to (0, 1)$
- $g'(z) = g(z)(1 - g(z))$
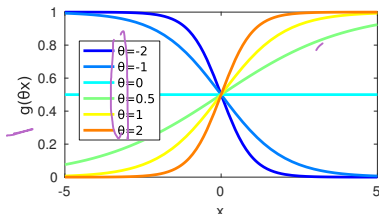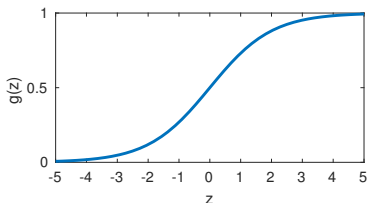
# Logistic Regression Hypothesis Function

### Sigmoid function

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

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



Hypothesis function for logistic regression:

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

$$\theta^T x = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n-1} \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix}$$

# Review: Bernoulli Distribution

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

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

$x = 1, \ p(x) = \lambda (1-\lambda)^0 = \lambda$

$x = 0, \ p(x) = \lambda^0 (1-\lambda)^1 = 1 - \lambda$

landing on $(H, H)$

$0.7 \cdot 0.7 = 0.49$

$(H, T): \ 0.7 \cdot (1-0.7) = 0.21$

$\lambda = 0.7$

$P(\text{head}) = P(x=1) = \lambda = 0.7$

$P(\text{tail}) = P(x=0) = 1 - \lambda = 0.3$

# Maximum likelihood estimation for logistic regression

*model for* $\lambda$     $\varepsilon^{(i)}$       $\mathbb{E}[y|x] = \lambda$
$= h_\theta(x)$

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

- $p(y = 1 \mid x; \theta) = h_\theta(x)$
- $p(y = 0 \mid x; \theta) = 1 - h_\theta(x)$

*model for* $1 - \lambda$

# Maximum likelihood estimation for logistic regression

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

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

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

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

# Maximum likelihood estimation for logistic regression

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

- $p(y = 1 \mid x; \theta) = h_\theta(x)$
- $p(y = 0 \mid x; \theta) = 1 - h_\theta(x)$

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

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

*i.i.d.*

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

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

# Maximum likelihood estimation for logistic regression

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

- $p(y = 1 \mid x; \theta) = h_\theta(x)$
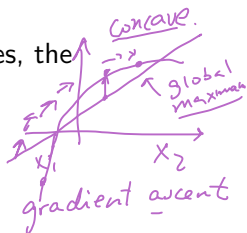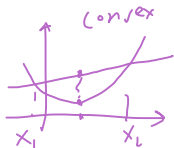- $p(y = 0 \mid x; \theta) = 1 - h_\theta(x)$

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

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!

*Handwritten annotations:* Convex $x_1$ $x_2$. Concave, global maximum $x_1$ $x_2$, gradient ascent.

# 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 $\text{argmax}_\theta\, l(\theta)$ using gradient ascent:

$$\frac{\partial l(\theta)}{\partial \theta_j} =$$

# 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 $\text{argmax}_\theta \, l(\theta)$ using gradient ascent:

*similar to* $\nabla J(\theta)$ *in linear regression* $h_\theta(x) = \theta^T x$

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

*sigmoid* $h_\theta(x^i) = \dfrac{1}{1 + e^{-\theta^T x}}$

## Stocastic Gradient Ascent

```
Repeat until convergence{
  for i = 1...m {
    θj = θj + α(y(i) − hθ(x(i)))x(i)j   for every j
  }
}
```

▶ Update rule has the same form as least square regression, but with different hypothesis function $h_\theta$

# Binary Digit Classification

$$0 < g(\theta^T x) < 1$$

### Using the learned classifier

Given an image $x$, the predicted label is

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

### Binary digit classification results

|  | sample size | accuracy |
|---|---|---|
| Training | 16200 | 100% |
| Testing | 1225 | 100% |

▶ Testing accuracy is 100% since this problem is relatively easy.
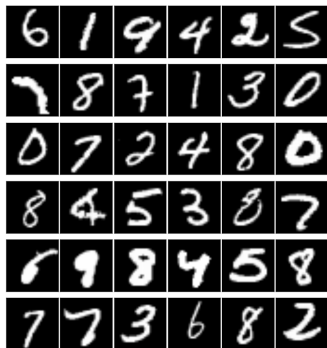
## Multi-Class Classification

# Multi-class classification

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

$$\mathcal{Y} = \{1, \ldots, k\} \qquad k = 10.$$

MNIST Samples



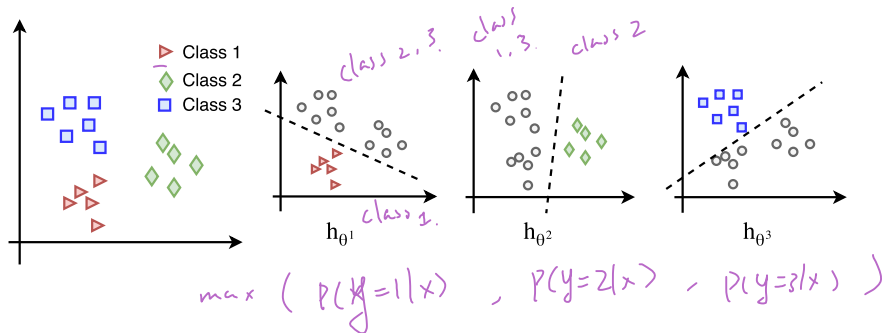Given new sample $x \in \mathbb{R}^k$, predict which class it belongs.

# Naive Approach: Convert to binary classification

## One-Vs-Rest

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

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

$$\hat{y} = \underset{i}{\arg\max}\, h_i(x)$$

# Multiple binary classifiers

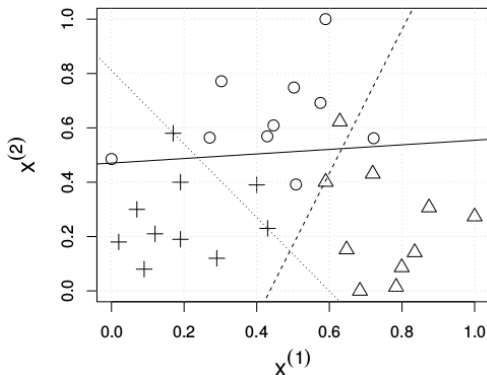50-class classification

$\left(\frac{1}{50}, \frac{49}{50}\right)$

Drawbacks of One-Vs-Rest:

▶ 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
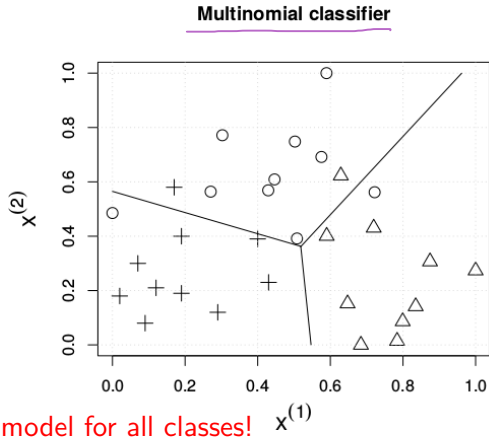- ▶ Different classifiers may have different confidence scales

**Multinomial classifier**



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 $\phi_i$

$$\sum_{i=1}^{k} \phi_i = 1$$

fair die

$$\phi = \left[ \tfrac{1}{6}, \tfrac{1}{6}, \tfrac{1}{6}, \tfrac{1}{6}, \tfrac{1}{6}, \tfrac{1}{6} \right]$$

$\phi$:

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

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

# Extend logistic regression: Softmax Regression

i.i.d.

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

# Extend logistic regression: Softmax Regression

$$h_\theta(x) = \begin{bmatrix} \lambda \\ 1-\lambda \end{bmatrix}$$

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

Hypothesis function for sample $x$:

$$p(y = \ell \mid x; \theta) = \frac{e^{\theta_\ell^T x}}{\sum_{\nu=1}^{k} e^{\theta_\nu^T x_j}}$$

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

$$z$$

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

# Extend logistic regression: Softmax Regression

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

Hypothesis function for sample $x$:

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

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

$\theta_{1,1} \cdots \quad \theta_{1,n}$

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

$\ominus$

$(k \times n)$

# Softmax Regression

Bernoulli $\quad p(y|x) = \underline{\lambda}^{y}(1-\lambda)^{\underline{y}}.$

$\lambda = h_\theta(x)$

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

$$1\{y^i = \ell\} = \begin{cases} 1 & y^i = \ell \\ 0 & y^i \neq \ell \end{cases}$$

$$\ell(\theta) = \sum_{i=1}^{m} \log \underline{p(y^{(i)}|x^{(i)}; \theta)}$$

$$= \sum_{i=1}^{m} \log \prod_{l=1}^{k} p(y^{(i)} = l|x^{(i)})^{\mathbf{1}\{y^{(i)} = l\}}$$

# Softmax Regression

Given $(x^{(i)}, y^{(i)}), i = 1, \ldots, m$, 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 $(x^{(i)}, y^{(i)}), i = 1, \ldots, m$, 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 $\nabla_{\theta_l} \ell(\theta)$

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

# Property of Softmax Regression

▶ Parameters $\theta_1, \ldots \theta_k$ are not independent:
$\sum_j p(y = j|x) = \sum_j \phi_j = 1$

▶ Knowning $k - 1$ 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 $\theta* = \theta - \begin{bmatrix} \theta_2 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \theta_1 - \theta_2 \\ 0 \end{bmatrix}$,

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

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

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

# When to use Softmax?

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