

CSC 4020 Fundamental of Machine Learning: Linear Regression

Baoyuan Wu
School of Data Science, CUHK-SZ

January 25/27, 2021

1 Linear Regression: A Deterministic Perspective

2 Linear Regression: A Probabilistic Perspective

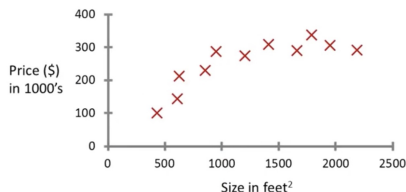
- Probabilistic modeling
- Robust linear regression
- Ridge regression
- Lasso regression

3 Generalized Linear Regression

Linear regression with one variable

- Here we start from a simple example of one dimensional input variable, and the training dataset $D = \{(x_i, y_i)\}_{i=1}^m$ can be plotted on the $x - y$ plane.
- m indicates the number of training samples; x denotes the input variable/feature; y denotes the output variable.

Size in feet ² (x)	Price in 1000's (y)
2104	460
1416	232
1514	315
852	178
...	...



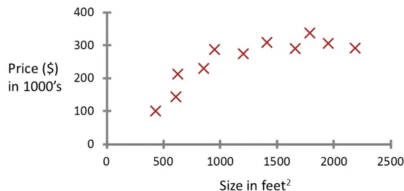
Linear hypothesis function

- Our goal is to find a linear hypothesis function to well fit the training data D , *i.e.*,

$$h_{\theta}(x) = \theta_0 + \theta_1 \phi(x) = [\theta_0, \theta_1][1; \phi(x)] = \hat{\phi}(x)^{\top} \theta \quad (1)$$

where $\phi(x)$ is called **basis expansion**, which is specified as different forms, such as $\phi(x) = x$ or $\phi(x) = [x^3; x^2; x]$. In the following, we will use $\phi(x) = x$ as example, while other expansions will be introduced later.

- Given θ_0, θ_1 , $h_{\theta}(x)$ is the function of x .
- Given x , $h_{\theta}(x)$ is a **linear function** of $\theta = [\theta_0; \theta_1]$. This is why it is called **linear regression**.
- Then, given D , how to learn θ ?



Cost function

- We design the following **cost function** to minimize the difference between the prediction $h_{\boldsymbol{\theta}}(x_i)$ and the ground-truth value y_i , *i.e.*,

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(x_i) - y_i)^2 \quad (2)$$

$$= \frac{1}{2} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2, \quad (3)$$

$$= \frac{1}{2} \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i)^2 \quad (4)$$

which is called **residual sum of squares** (RSS) or sum of squared errors (SSE).

- $J(\boldsymbol{\theta})$ is a **convex** or **non-convex** function? What is the shape of it?

Gradient descent

- The linear regression is formulated to the following optimization problem

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i)^2. \quad (5)$$

- $\boldsymbol{\theta}$ can be updated by **gradient descent algorithm**,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad \frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i) \bar{\mathbf{x}}_i \quad (6)$$

where α is called step-size or learning rate.

- Does gradient descent always converge to the optimal solution? (Plot the trajectory of gradient descent on curve or contours)

Analytical solution

- If we set the gradient to 0, then we can get the following solution

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i) \bar{\mathbf{x}}_i = \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^\top \mathbf{y} = 0 \quad (7)$$

$$\Rightarrow \boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (8)$$

which are called **normal equation** and **ordinary least squares** (OLS) solution, respectively. $\mathbf{X} = [\bar{\mathbf{x}}_1^\top; \bar{\mathbf{x}}_2^\top; \dots; \bar{\mathbf{x}}_m^\top] \in \mathbb{R}^{m \times d}$.

- Since there is a closed-form solution, why do we need gradient descent algorithm?

Geometric interpretation

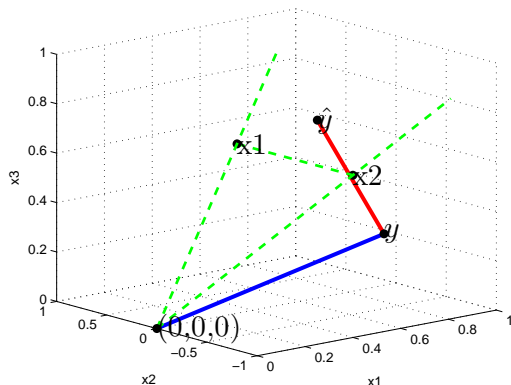
- Since $\theta^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$, then the predictions of \mathbf{X} can be obtained by

$$\hat{\mathbf{y}} = \mathbf{X}\theta^* = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (9)$$

which corresponds to the **orthogonal projection** of \mathbf{y} onto the column space of \mathbf{X} .

$$\mathbf{X} = \begin{pmatrix} 1 & 2 \\ 1 & -2 \\ 1 & 2 \end{pmatrix},$$

$$\mathbf{y} = \begin{pmatrix} 8.89 \\ 0.61 \\ 1.77 \end{pmatrix}$$



Normal equation vs. gradient descent

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

m -dimensional vector

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

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$, need to calculate inverse of $X^T X$
Works well when n is large	Slow if n is very large

Probabilistic modeling

- We assume that the relationship between the input variable/feature \mathbf{x} and the output variable y is

$$y = \boldsymbol{\theta}^\top \mathbf{x} + e, \text{ where } e \sim \mathcal{N}(0, \sigma^2), \quad (10)$$

where e is called **observation noise** or **residual error**, and it is independent with any specific input \mathbf{x} .

- Thus, the output y can also be seen as a random variable, and its conditional probability is formulated as

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}, \sigma^2) \quad (11)$$

Maximum log-likelihood estimation

- The parameter θ can be learned by maximum log-likelihood estimation (MLE), given the training dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, as follows

$$\theta_{MLE} = \arg \max_{\theta} \log \mathcal{L}(\theta|D) \quad (12)$$

$$= \sum_i^m \log p(y|\mathbf{x}, \theta) = \sum_i^m \log \mathcal{N}(\theta^\top \mathbf{x}, \sigma^2) \quad (13)$$

$$= -\log(\sigma^m (2\pi)^{\frac{m}{2}}) - \frac{1}{2\sigma^2} \sum_i^m (y_i - \theta^\top \mathbf{x}_i)^2 \quad (14)$$

- Removing the constants w.r.t. θ ,

$$\theta_{MLE} = \arg \min_{\theta} \frac{1}{2} \sum_i^m (y_i - \theta^\top \mathbf{x}_i)^2, \quad (15)$$

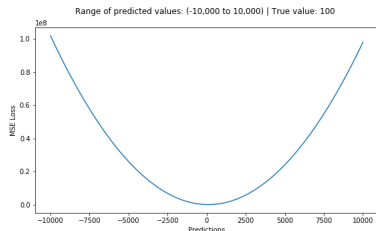
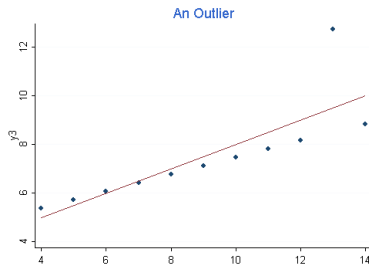
which is exactly same with the cost function from the deterministic perspective.

Robust linear regression

- When there is a few outliers in the training data D , which are far from most other points, then learned parameters θ_{MLE} will be significantly influenced, leading to very poor fit.
- Let's see the loss curve of the residual sum of squares (RSS),

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\bar{x}_i^\top \theta - y_i)^2. \quad (16)$$

- The error increases quadratically along with the residual. To minimize such a large error, the linear model will be significantly changed.
- How to alleviate the significant influence of outliers?

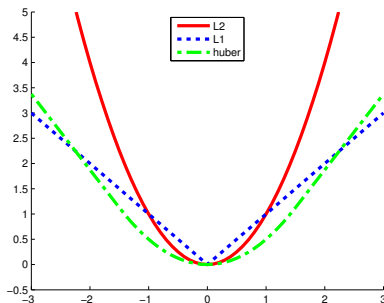


Robust linear regression

- We adopt the ℓ_1 loss to replace the ℓ_2 loss, as follows

$$J(\boldsymbol{\theta}) = \sum_{i=1}^m |\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i|. \quad (17)$$

- The curves of ℓ_1 and ℓ_2 losses are shown ad follows.
- When the residual is large, the ℓ_1 loss is much smaller than the ℓ_2 loss, such that the influence of outliers could be alleviated.



Robust linear regression

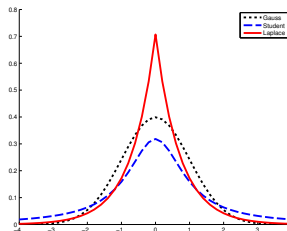
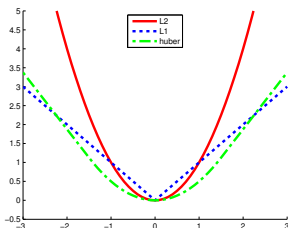
- Actually, the above ℓ_1 loss can also be derived from the probabilistic perspective, by assuming that

$$p(y|\mathbf{x}, \boldsymbol{\theta}, b) = \text{Lap}(y|\mathbf{x}, \boldsymbol{\theta}, b) \propto \exp\left(-\frac{1}{b}|y - \boldsymbol{\theta}^\top \mathbf{x}|\right) \quad (18)$$

- Applying the maximum log-likelihood estimation (MLE), we will obtain

$$\boldsymbol{\theta}_{MLE} = \arg \max_{\boldsymbol{\theta}} \log \mathcal{L}(\boldsymbol{\theta}|D) = \sum_i^m \log p(y_i|\mathbf{x}_i, \boldsymbol{\theta}) \quad (19)$$

$$\equiv \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m |\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i| \quad (20)$$



Robust linear regression

$$\boldsymbol{\theta}_{MLE} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m |\mathbf{x}_i^{\top} \boldsymbol{\theta} - y_i| \quad (21)$$

- However, the ℓ_1 loss function is non-differentiable and non-linear. The gradient descent algorithm cannot be adopted.
- We can transform it to a linear program, as follows

$$\min_{\boldsymbol{\theta}, \mathbf{t}} \sum_i^m t_i \quad (22)$$

$$s.t. \quad -t_i \leq \mathbf{x}_i^{\top} \boldsymbol{\theta} - y_i \leq t_i, 1 \leq i \leq m. \quad (23)$$

Robust linear regression

$$\boldsymbol{\theta}_{MLE} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m |\mathbf{x}_i^\top \boldsymbol{\theta} - y_i| \quad (24)$$

- We can also utilize the following equation:

$$|a| = \min_{\mu} \frac{1}{2} \left(\frac{a^2}{\mu} + \mu \right) \quad (25)$$

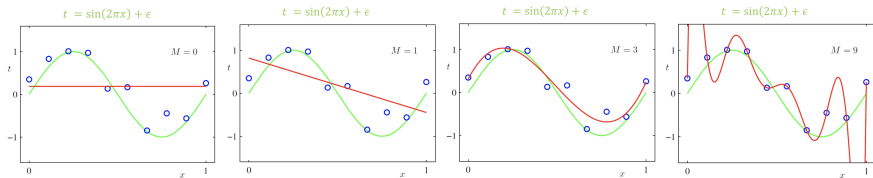
- Then, the ℓ_1 minimization problem can be reformulated as follows

$$\min_{\boldsymbol{\theta}} \min_{\mu} \frac{1}{2} \left(\frac{(\mathbf{x}^\top \boldsymbol{\theta} - y_i)^2}{\mu} + \mu \right). \quad (26)$$

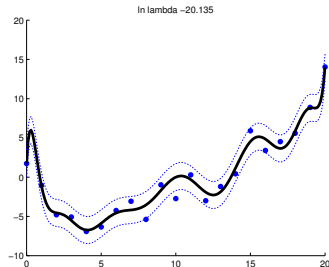
- It can be iteratively and alternatively optimized as follows:
 - Given $\boldsymbol{\theta}$, $\mu = |\mathbf{x}^\top \boldsymbol{\theta} - y_i|$
 - Given μ , $\boldsymbol{\theta} = \min_{\boldsymbol{\theta}} (\mathbf{x}^\top \boldsymbol{\theta} - y_i)^2$
- It is called **iteratively reweighted least squares** method.

Ridge regression

- As demonstrated in the first week, overfitting is an important challenge for linear regression.
- What approaches we have introduced to alleviate overfitting? Ocam's razor or cross-validation
- Is there other more theoretical approaches? SURE!



Ridge regression



- Let's see one simple example, we use a polynomial function with 14 degree to fit $m = 21$ data points. The learned curve is very “wiggly” (see above).
- The parameter values of this curve are as follows

6.56, -36.934, -109.25, 543.452, 1022.561, -3046.224, -3768.013, 8524.54, 6607.897, -12640.058, -5530.188, 9479.73, 1774, 639, -2821.526

- There are many large positive/negative values, such that a small change of features could lead to significant change of output.

Ridge regression

- How to get smaller parameter values?
- We can assume that the parameter follow a zero-mean Gaussian prior

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \tau^2\mathbf{I}) \quad (27)$$

- Utilizing this prior, we obtain the maximum a posteriori (MAP) estimation

$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} \sum_i^m \log p(y|\mathbf{x}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (28)$$

$$= \sum_i^m \log \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}, \sigma^2) + \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \tau^2\mathbf{I}) \quad (29)$$

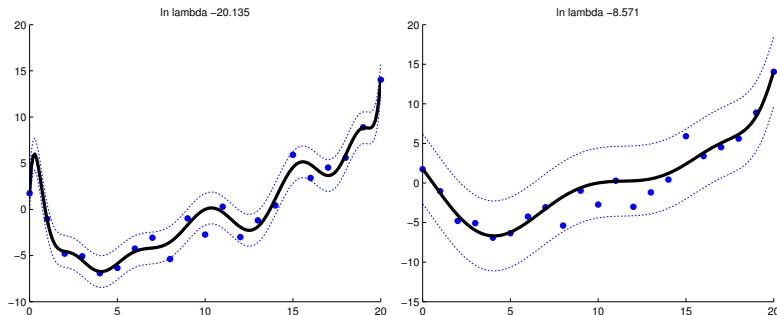
$$\equiv \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i)^2 + \lambda \|\boldsymbol{\theta}\|_2^2. \quad (30)$$

- The corresponding closed-form solution is given by

$$\boldsymbol{\theta}_{MAP} = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (31)$$

Ridge regression

- The above method is also known as **ridge regression**, or **penalized least squares**.
- In general, adding a Gaussian prior to the parameters of a model to encourage them to be small is called ℓ_2 **regularization** or **weight decay**.
- As shown below, when we set a larger λ , *i.e.*, more weight on the prior, the resulting curve will be smoother.



Lasso regression

- We can replace the Gaussian prior by a Laplacian prior, *i.e.*,

$$p(\boldsymbol{\theta}) = \text{Lap}(\boldsymbol{\theta}|\mathbf{0}, b) = \frac{1}{2b} \exp\left(-\frac{|\boldsymbol{\theta}|}{b}\right), \quad (32)$$

- The combination of the Gaussian distribution of $p(y|\mathbf{x}, \boldsymbol{\theta})$ and the Laplacian prior, leading to

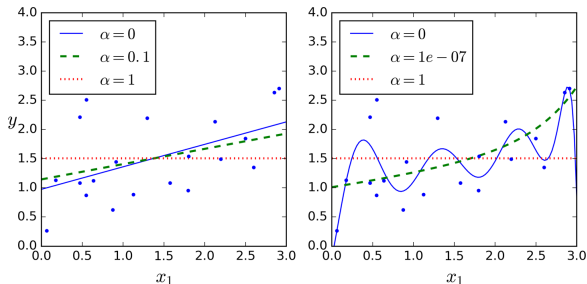
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} \sum_i^m \log p(y|\mathbf{x}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (33)$$

$$= \sum_i^m \log \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}, \sigma^2) + \text{Lap}(\boldsymbol{\theta}|\mathbf{0}, b) \quad (34)$$

$$\equiv \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m (\bar{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i)^2 + \lambda |\boldsymbol{\theta}|. \quad (35)$$

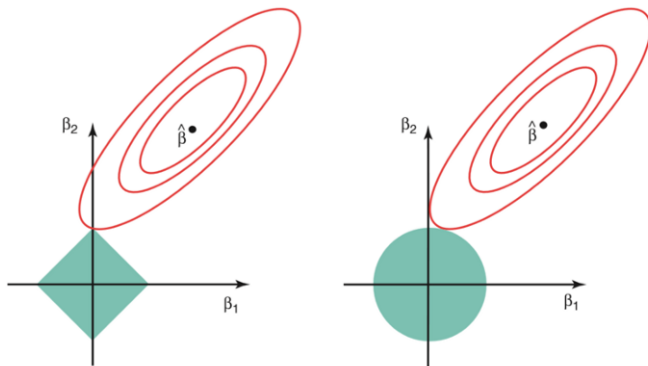
Lasso regression

- It is **Lasso regression**, and the regularization is called ℓ_1 **regularization**. It will encourage the sparse parameters.
- As shown below, when we set a larger λ , *i.e.*, more weight on the prior, the resulting curve will be smoother.



Geometry of Ridge and Lasso regression

- Geometry of Ridge and Lasso regression. Which one is Ridge?

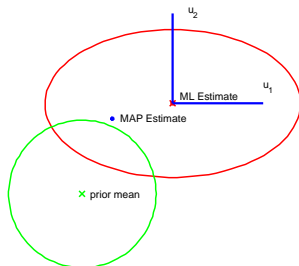


Summary of different linear regressions

Note that the uniform distribution will not change the mode of the likelihood.

Thus, MAP estimation with a uniform prior corresponds to MLE.

$p(y \mathbf{x}, \boldsymbol{\theta})$	$p(\boldsymbol{\theta})$	regression method
Gaussian	Uniform	Least squares
Gaussian	Gaussian	Ridge regression
Gaussian	Laplace	Lasso regression
Laplace	Uniform	Robust regression
Student	Uniform	Robust regression



Generalized linear regression

- **Linear model:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \phi(\mathbf{x}), \quad (36)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (37)$$

where f denotes a distribution function.

- **Generalized linear model (GLM):**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad (38)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (39)$$

where g is called **link function**, which is required to be monotonically increasing differentiable.

- The standard linear model is a special case of GLM with $g(a) = a$.

Why we need generalized linear regression

- **Why we need generalized linear model?** Let's see one example.

In the early stages of a disease epidemic, the rate at which new cases occur can often increase exponentially through time. Hence, if μ_i is the expected number of new cases on day t_i , a model of the form

$$\mu_i = \gamma \exp(\delta t_i)$$

seems appropriate.

- ▶ Such a model can be turned into GLM form, by using a **log link** so that

$$\log(\mu_i) = \log(\gamma) + \delta t_i = \beta_0 + \beta_1 t_i.$$

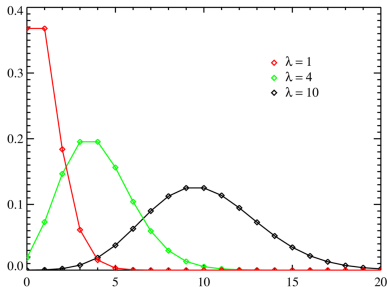
- ▶ Since this is a count, the **Poisson distribution** (with expected value μ_i) is probably a reasonable distribution to try.

Log linear regression

- **Poisson distribution** The Poisson distribution is popular for modeling the number of times an event occurs in an interval of time or space.
- A discrete random variable X is said to have a Poisson distribution with parameter $\lambda > 0$ if for $k = 0, 1, 2, \dots$, the probability mass function of X is given by

$$f(k; \lambda) = P(X = k | \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad (40)$$

where e is Euler's number ($e = 2.71828\dots$), k is the number of occurrences, $k!$ is the factorial of k .



Log linear regression

- We assume that the conditional probability follows

$$P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = \text{Poisson}(\lambda_i) = \frac{\lambda^{y_i} e^{-\lambda}}{y_i!}, \quad \ln \lambda_i = \boldsymbol{\theta}^\top \mathbf{x}_i \quad (41)$$

- The log-likelihood function is formulated as follows

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^m \log P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = \sum_{i=1}^m y_i \log \lambda_i - \lambda_i - \log y_i! \quad (42)$$

- We have

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (y_i \mathbf{x}_i - e^{\boldsymbol{\theta}^\top \mathbf{x}_i}) = 0 \quad \Rightarrow \quad \ln y_i = (\boldsymbol{\theta}^*)^\top \mathbf{x}_i \quad (43)$$

- Plot the log-linear regression as below.

Logistic regression

- We assume that the conditional probability follows

$$P(y_i|\mathbf{x}_i, \boldsymbol{\theta}, N) = \text{Bin}(y_i|N, \mu_i) = \binom{N}{y_i} \mu_i^{y_i} (1 - \mu_i)^{N - y_i}, \quad \mu_i = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}. \quad (44)$$

- The log-likelihood function is formulated as follows

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^m \log P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = y_i \log \mu_i + (N - y_i) \log(1 - \mu_i) \quad (45)$$

- We have

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (y_i - N\mu_i) \mathbf{x}_i = 0 \quad \Rightarrow \quad \frac{y_i}{N} = \mu_i = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}. \quad (46)$$

- Since the $\sigma(a) = \frac{1}{1+e^{-a}}$ is called **sigmoid function** or **logit function**, the above model is called **logit regression** or **logistic regression**.
- Since $\frac{y_i}{N} \in [0, 1]$, it can be seen as the posterior probability. Thus, logistic regression is a classification model, rather than regression.

- Linear model is the linear function of the parameter θ , rather than the input feature
- Linear model is a special case of generalized linear model, while generalized linear model is not always linear
- Choosing different linear models is equivalent to choosing different distributions of $p(y|\mathbf{x}, \theta)$ and $p(\theta)$, according to the task and the data

Reading material

- <https://www.stat.cmu.edu/~ryantibs/advmethods/notes/glm.pdf>