

CSC 4020 Fundamentals of Machine Learning: Linear Regression

Baoyuan Wu
School of Data Science, CUHK-SZ

February 22/24, 2021

Outline

- 1 Review of last week
- 2 Classification and representation
- 3 Logistic regression
- 4 Regularized logistic regression

Linear regression: deterministic perspective

Linear regression: deterministic perspective

- Linear hypothesis function: $h_{\theta} = \phi(\mathbf{x})^{\top} \theta$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$
$$\phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_1 x_2 \\ \vdots \end{bmatrix}$$

Linear regression: deterministic perspective

- **Linear hypothesis function:** $h_{\theta} = \phi(\mathbf{x})^{\top} \theta$
- **Linear regression** by minimizing residual sum of squares (RSS):

$$\theta^* = \arg \min_{\theta} J(\theta) = \frac{1}{2} \sum_{i=1}^m (\underbrace{\phi(\mathbf{x})_i^{\top} \theta}_{\text{}} - \underbrace{y_i}_{\text{}})^2$$

Linear regression: deterministic perspective

- **Linear hypothesis function:** $h_{\theta} = \phi(\mathbf{x})^{\top} \theta$
- **Linear regression** by minimizing residual sum of squares (RSS):

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

- **Two solutions:** gradient descent and close-form solution (called normal equation or ordinary least squares solution)

Handwritten red notes and equations illustrating the normal equation derivation:

- $(\phi(\mathbf{x})^{\top} \theta - y) \cdot \phi(\mathbf{x}) = 0$
- $\phi^{\top}(\mathbf{x}) \phi(\mathbf{x}) \theta = y \phi(\mathbf{x})$
- $\phi(\mathbf{x})_m \theta$
- $\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$

Linear regression: probabilistic perspective

Linear regression: probabilistic perspective

- We assume that: $y = \underbrace{\boldsymbol{\theta}^\top \mathbf{x}} + \underbrace{e}$, where $e \sim \underbrace{\mathcal{N}(0, \sigma^2)}$ is called **observation noise** or **residual error**

Linear regression: probabilistic perspective

- We assume that: $y = \boldsymbol{\theta}^\top \mathbf{x} + e$, where $e \sim \mathcal{N}(0, \sigma^2)$ is called **observation noise** or **residual error**
- y is also a random variable, and its conditional probability is

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

Linear regression: probabilistic perspective

- We assume that: $y = \boldsymbol{\theta}^\top \mathbf{x} + e$, where $e \sim \mathcal{N}(0, \sigma^2)$ is called **observation noise** or **residual error**
- y is also a random variable, and its conditional probability is

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

- **Maximum log-likelihood estimation:**

$$\boldsymbol{\theta}_{MLE} = \arg \max_{\boldsymbol{\theta}} \log \mathcal{L}(\boldsymbol{\theta}|D) \quad (1)$$

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

$$= -\log(\sigma^m (2\pi)^{\frac{m}{2}}) - \frac{1}{2\sigma^2} \sum_i^m (y_i - \boldsymbol{\theta}^\top \mathbf{x}_i)^2 \quad (3)$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \sum_i^m (y_i - \boldsymbol{\theta}^\top \mathbf{x}_i)^2, \quad (4)$$

Variants of linear regression

- Robust regression for data with outliers: $\theta_{MLE} = \arg \min_{\theta} \sum_{i=1}^m |\bar{x}_i^{\top} \theta - y_i|$
- Ridge regression to avoid over-fitting, through MAP estimation:

$$\theta_{MAP} = \arg \max_{\theta} \sum_{i=1}^m \log p(y|\mathbf{x}, \theta) + \log p(\theta) \quad (5)$$

$$= \sum_{i=1}^m \log \mathcal{N}(\theta^{\top} \mathbf{x}, \sigma^2) + \mathcal{N}(\theta | \mathbf{0}, \tau^2 \mathbf{I}) \quad (6)$$

$$\equiv \arg \min_{\theta} \sum_{i=1}^m (\bar{x}_i^{\top} \theta - y_i)^2 + \lambda \|\theta\|_2^2. \quad (7)$$

- Lasso regression to obtain sparse model,

$$\theta_{MAP} = \arg \max_{\theta} \sum_{i=1}^m \log \mathcal{N}(\theta^{\top} \mathbf{x}, \sigma^2) + \text{Lap}(\theta | \mathbf{0}, b) \quad (8)$$

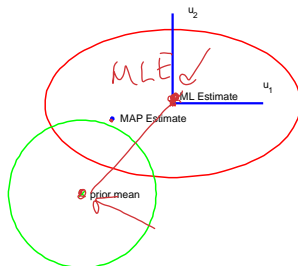
$$= \arg \min_{\theta} \sum_{i=1}^m (\bar{x}_i^{\top} \theta - y_i)^2 + \lambda |\theta|. \quad (9)$$

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.

$p(y x, \theta)$	$p(\theta)$	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

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \underbrace{g^{-1}}_{\triangle}(\underbrace{\boldsymbol{\theta}^\top \phi(\mathbf{x})}), \quad y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (10)$$

where g is called link function.

Generalized linear regression

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (10)$$

where g is called **link function**.

- We assume that the conditional probability follows

$$P(\underline{y_i}|\underline{x_i}, \boldsymbol{\theta}, N) = \underline{\text{Bin}}(\underline{y_i}|N, \mu_i) = \binom{N}{y_i} \mu_i^{y_i} (1 - \mu_i)^{N-y_i}, \quad (11)$$

- Where $\mu_i = \frac{1}{1+e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}, g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i}$.

Generalized linear regression

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (10)$$

where g is called **link function**.

- We assume that the conditional probability follows

$$P(y_i|\mathbf{x}_i, \boldsymbol{\theta}, N) = \text{Bin}(y_i|N, \mu_i) = \binom{N}{y_i} \mu_i^{y_i} (1 - \mu_i)^{N-y_i}, \quad (11)$$

- Where $\mu_i = \frac{1}{1+e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}, g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i}$.
- The log-likelihood function is formulated as follows

$$\log \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^m \log P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = y_i \log \mu_i + (N - y_i) \log(1 - \mu_i) \quad (12)$$

Generalized linear regression

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (10)$$

where g is called **link function**.

- We assume that the conditional probability follows

$$P(y_i|\mathbf{x}_i, \boldsymbol{\theta}, N) = \text{Bin}(y_i|N, \mu_i) = \binom{N}{y_i} \mu_i^{y_i} (1 - \mu_i)^{N-y_i}, \quad (11)$$

- Where $\mu_i = \frac{1}{1+e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}, g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i}$.
- The log-likelihood function is formulated as follows

$$\log \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^m \log P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = y_i \log \mu_i + (N - y_i) \log(1 - \mu_i) \quad (12)$$

- We have

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (y_i - N\mu_i) \mathbf{x}_i = 0 \Rightarrow \frac{y_i}{N} = \mu_i = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}. \quad (13)$$

Generalized linear regression

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (10)$$

where g is called **link function**.

- We assume that the conditional probability follows

$$P(y_i|\mathbf{x}_i, \boldsymbol{\theta}, N) = \text{Bin}(y_i|N, \mu_i) = \binom{N}{y_i} \mu_i^{y_i} (1 - \mu_i)^{N-y_i}, \quad (11)$$

- Where $\mu_i = \frac{1}{1+e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}, g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i}$.
- The log-likelihood function is formulated as follows

$$\log \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^m \log P(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = y_i \log \mu_i + (N - y_i) \log(1 - \mu_i) \quad (12)$$

- We have

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^m (y_i - N\mu_i) \mathbf{x}_i = 0 \quad \Rightarrow \quad \frac{y_i}{N} = \mu_i = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}_i}}. \quad (13)$$

- 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**.

Summary of last week

- Linear model is the linear function of the parameter θ , rather than the input feature

Summary of last week

- 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

Summary of last week

- 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 $\underbrace{p(y|\mathbf{x}, \theta)}$ and $\underbrace{p(\theta)}$, according to the task and the data

Classification

- Classification: classifying input data into discrete states

Classification

- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?

Classification

- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?
 - ~~Whether~~ forecast: sunny / not sunny?

Classification

- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam? $\rightarrow \{0, 1\}$
 - Whether forecast: sunny / not sunny?
 - Tumor: malignant / not malignant?

Classification

- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?
 - Whether forecast: sunny / not sunny?
 - Tumor: malignant / not malignant?
- The label $y \in \{0, 1\}$:

Classification

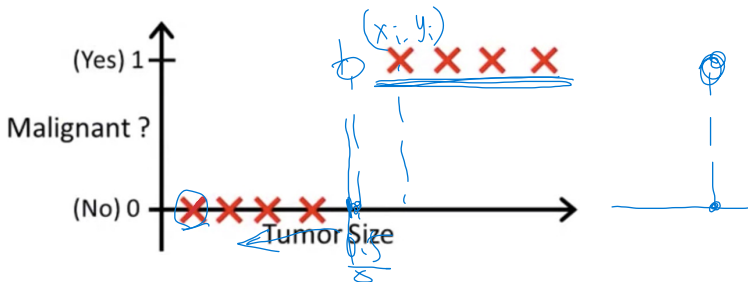
- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?
 - Whether forecast: sunny / not sunny?
 - Tumor: malignant / not malignant?
- The label $y \in \{0, 1\}$:
 - $y = 0$: negative class, *e.g.*, not spam, not sunny, not malignant

Classification

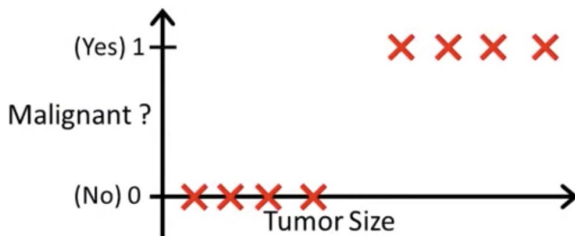
- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?
 - Whether forecast: sunny / not sunny?
 - Tumor: malignant / not malignant?
- The label $y \in \{0, 1\}$:
 - $y = 0$: negative class, *e.g.*, not spam, not sunny, not malignant
 - $y = 1$: positive class, *e.g.*, spam, sunny, malignant

Classification

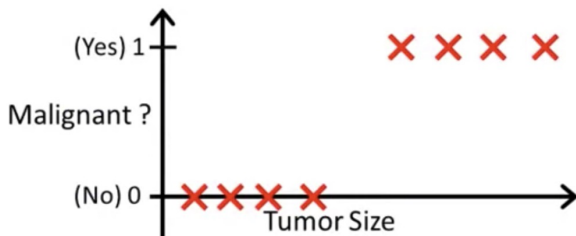
- Classification: classifying input data into discrete states
 - Email filtering: spam / not spam?
 - Weather forecast: sunny / not sunny?
 - Tumor: malignant / not malignant?
- The label $y \in \{0, 1\}$:
 - $y = 0$: negative class, *e.g.*, not spam, not sunny, not malignant
 - $y = 1$: positive class, *e.g.*, spam, sunny, malignant



Threshold classifier with linear regression

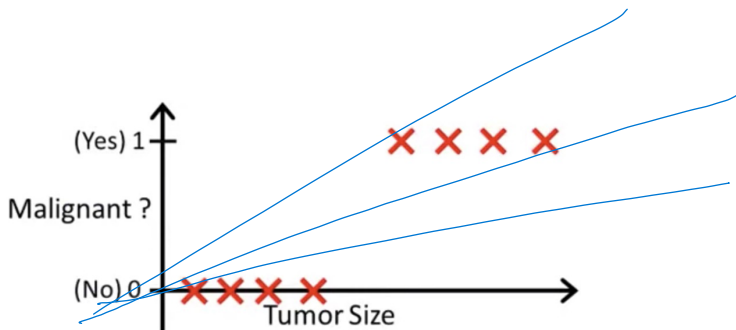


Threshold classifier with linear regression



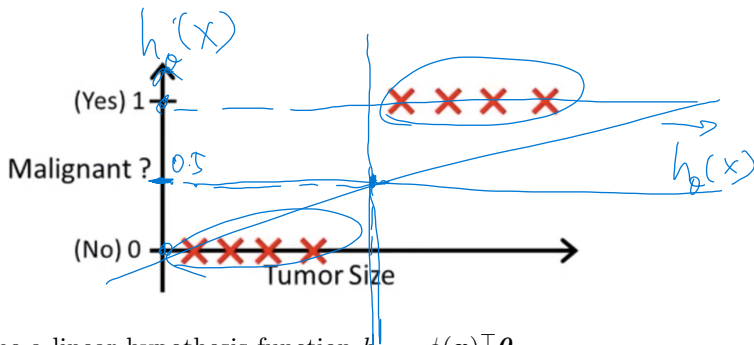
- We assume a linear hypothesis function $h_{\theta} = \underbrace{\phi(\mathbf{x})^T}_{\mathbf{x}} \boldsymbol{\theta}$

Threshold classifier with linear regression



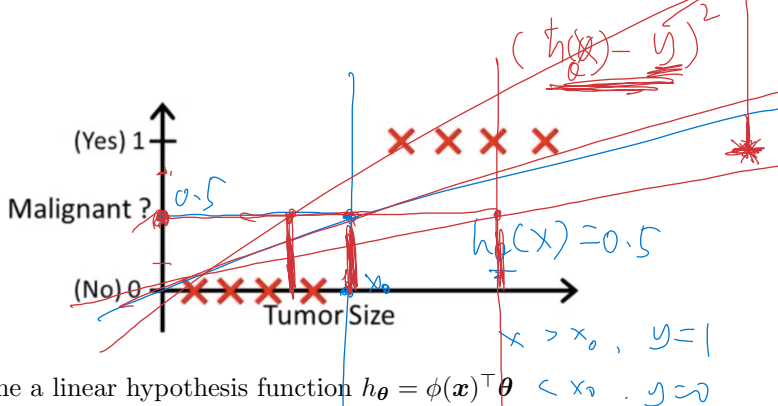
- We assume a linear hypothesis function $h_{\theta} = \phi(x)^{\top} \theta$
- A simple threshold classifier with this hypothesis function is

Threshold classifier with linear regression



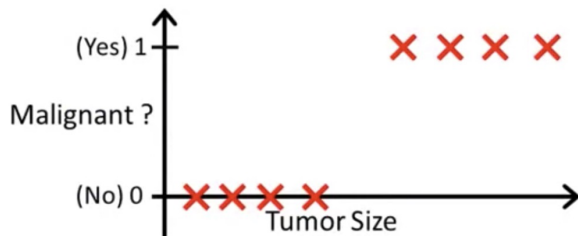
- We assume a linear hypothesis function $h_{\theta} = \phi(x)^{\top} \theta$
- A simple threshold classifier with this hypothesis function is
 - If $h_{\theta} > 0.5$, then $y = 1$, i.e., malignant tumor

Threshold classifier with linear regression



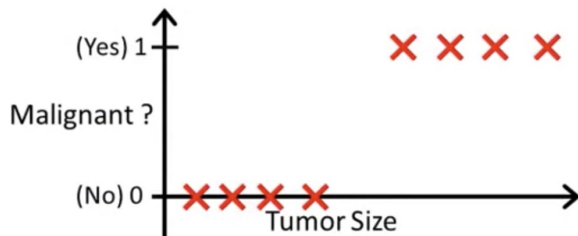
- We assume a linear hypothesis function $h_\theta = \phi(x)^\top \theta$
- A simple threshold classifier with this hypothesis function is
 - If $h_\theta > 0.5$, then $y = 1$, i.e., malignant tumor
 - If $h_\theta < 0.5$, then $y = 0$, i.e., benign tumor

Threshold classifier with linear regression



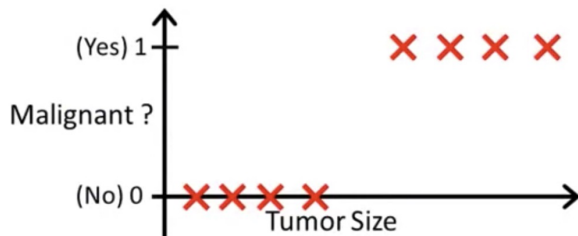
- It seems that the simple threshold classifier with linear regression works well on this classification task

Threshold classifier with linear regression



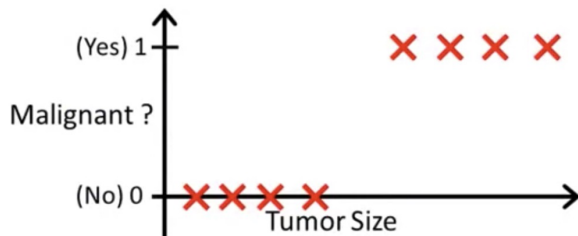
- It seems that the simple threshold classifier with linear regression works well on this classification task
- However, if there is a positive sample with very large tumor size (plot above), what will happen?

Threshold classifier with linear regression



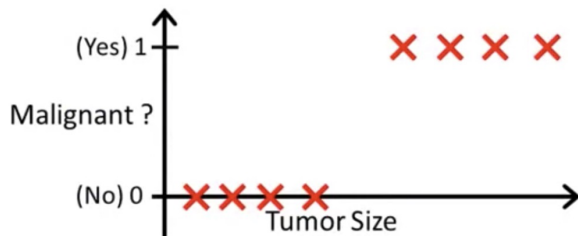
- It seems that the simple threshold classifier with linear regression works well on this classification task
- However, if there is a positive sample with very large tumor size (plot above), what will happen?
- The hypothesis function will be significantly changed, causing that some positive samples are mis-classified as negative (not malignant)?

Threshold classifier with linear regression



- It seems that the simple threshold classifier with linear regression works well on this classification task
- However, if there is a positive sample with very large tumor size (plot above), what will happen?
- The hypothesis function will be significantly changed, causing that some positive samples are mis-classified as negative (not malignant)? How to handle it?

Threshold classifier with linear regression



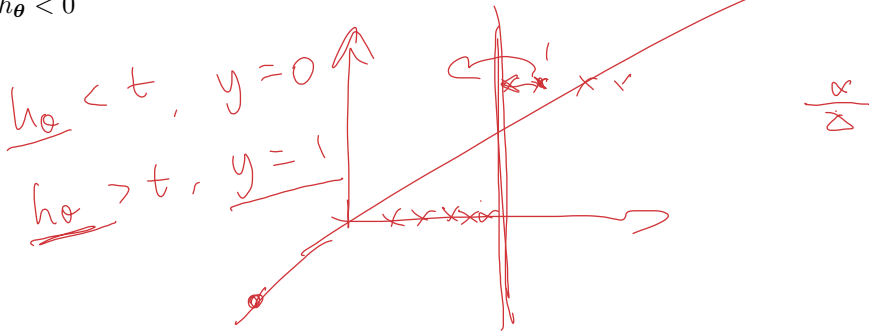
- It seems that the simple threshold classifier with linear regression works well on this classification task
- However, if there is a positive sample with very large tumor size (plot above), what will happen?
- The hypothesis function will be significantly changed, causing that some positive samples are mis-classified as negative (not malignant)? How to handle it? Adjusting the threshold value, or adopting robust linear regression.

Threshold classifier with linear regression

- But there is still something wired.

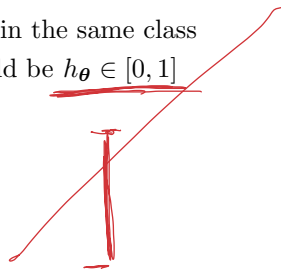
Threshold classifier with linear regression

- But there is still something wired.
- Our goal is to predict $y \in \{0, 1\}$, but the prediction could be $h_{\theta} > 1$ or $h_{\theta} < 0$



Threshold classifier with linear regression

- But there is still something wired.
- Our goal is to predict $y \in \{0, 1\}$, but the prediction could be $h_{\theta} > 1$ or $h_{\theta} < 0$
- It cannot reflect the difference among samples within the same class
- An expected hypothesis function for this task should be $h_{\theta} \in [0, 1]$



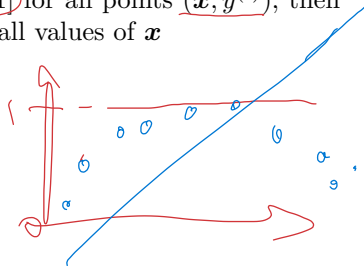
Threshold classifier with linear regression

Which statements are true?

Threshold classifier with linear regression

Which statements are true?

- If linear regression doesn't work well like the above example, feature scaling may help
- If the training set satisfies that all $y^{(i)} \in [0, 1]$ for all points $(\mathbf{x}, y^{(i)})$, then the linear hypothesis function $h_{\theta} \in [0, 1]$ for all values of \mathbf{x}
- None of above two states are true



Hypothesis representation

- An expected hypothesis function for this task should be $h_{\theta} \in [0, 1]$

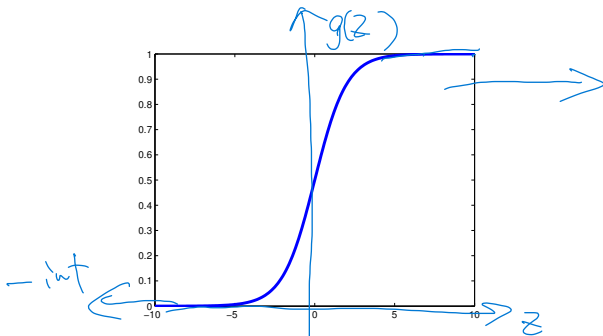
Hypothesis representation

- An expected hypothesis function for this task should be $h_{\theta} \in [0, 1]$
- Recall the generalized linear regression that

$$h_{\theta}(x) = \underbrace{g}_{\text{logit}}(\underbrace{\theta^T x}_u) \in [0, 1], \quad g(z) = \frac{1}{1 + \exp(-z)}$$

Handwritten notes: "logit" under g, "u" under theta^T x, "0" with an arrow pointing to the denominator of g(z).

where $g(\cdot)$ is called **sigmoid function** or **logistic regression**. (Plot below)



Hypothesis representation

- Interpretation of logistic function:

Hypothesis representation

- Interpretation of logistic function:
- $h_{\theta}(\mathbf{x}) =$ estimated probability that $y = 1$ of input \mathbf{x}

Hypothesis representation

- Interpretation of logistic function:
- $h_{\theta}(\mathbf{x})$ = estimated probability that $y = 1$ of input \mathbf{x}
- For example (plot below), if $h_{\theta}(\underline{\mathbf{x}}) = \underline{0.8}$, then it means that an patient with tumor size x has 80% chance of tumor being malignant

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e(-\theta^T \mathbf{x})} = \frac{p(y=1 | \mathbf{x}; \theta)}{1}$$

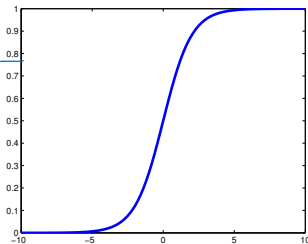
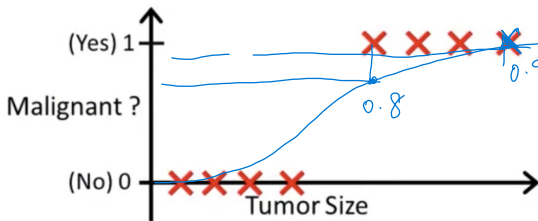
Hypothesis representation

- Interpretation of logistic function:
- $h_{\theta}(\mathbf{x})$ = estimated probability that $y = 1$ of input \mathbf{x}
- For example (plot below), if $h_{\theta}(\mathbf{x}) = 0.8$, then it means that an patient with tumor size x has 80% chance of tumor being malignant
- In this task, larger tumor size has larger chance/probability of being malignant tumor.

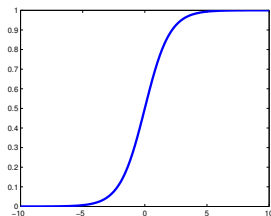
Hypothesis representation

- Interpretation of logistic function:
- $h_{\theta}(\mathbf{x})$ = estimated probability that $y = 1$ of input \mathbf{x}
- For example (plot below), if $h_{\theta}(\mathbf{x}) = 0.8$, then it means that a patient with tumor size x has 80% chance of tumor being malignant
- In this task, larger tumor size has larger chance/probability of being malignant tumor.
- Thus, we can say that

$$h_{\theta}(\mathbf{x}) = P(y = 1 | \mathbf{x}; \theta).$$



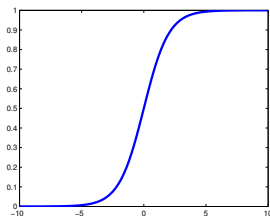
Decision boundary



- In logistic regression, we have

$$\begin{cases} h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^{\top} \mathbf{x}) = P(y = 1 | \mathbf{x}; \boldsymbol{\theta}) \in [0, 1], & (14) \\ g(z) = \frac{1}{1 + \exp(-z)}, & (15) \end{cases}$$

Decision boundary



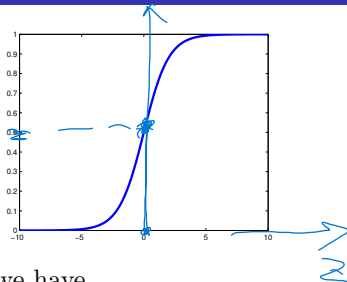
- In logistic regression, we have

$$h_{\theta}(\mathbf{x}) = g(\theta^{\top} \mathbf{x}) = P(y = 1 | \mathbf{x}; \theta) \in [0, 1], \quad (14)$$

$$g(z) = \frac{1}{1 + \exp(-z)}, \quad (15)$$

- Suppose that if $h_{\theta}(\mathbf{x}) > 0.5$, then we predict $y = 1$; if $h_{\theta}(\mathbf{x}) < 0.5$, then we predict $y = 0$

Decision boundary



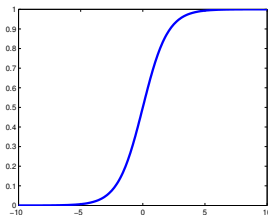
- In logistic regression, we have

$$h_{\theta}(\mathbf{x}) = g(\theta^{\top} \mathbf{x}) = P(y = 1 | \mathbf{x}; \theta) \in [0, 1], \quad (14)$$

$$g(z) = \frac{1}{1 + \exp(-z)}, \quad \geq 0.5 \quad (15)$$

- Suppose that if $h_{\theta}(\mathbf{x}) \geq 0.5$, then we predict $y = 1$; if $h_{\theta}(\mathbf{x}) < 0.5$, then we predict $y = 0$
- Correspondingly, if $\theta^{\top} \mathbf{x} \geq 0$, we predict $y = 1$; if $\theta^{\top} \mathbf{x} < 0$, then we predict $y = 0$.

Decision boundary

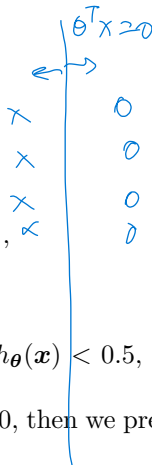


- In logistic regression, we have

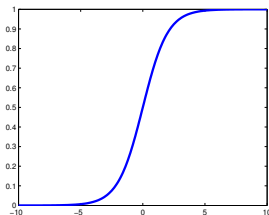
$$h_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x}) = P(y = 1 | \mathbf{x}; \theta) \in [0, 1], \quad (14)$$

$$g(z) = \frac{1}{1 + \exp(-z)}, \quad (15)$$

- Suppose that if $h_{\theta}(\mathbf{x}) \geq 0.5$, then we predict $y = 1$; if $h_{\theta}(\mathbf{x}) < 0.5$, then we predict $y = 0$
- Correspondingly, if $\theta^T \mathbf{x} \geq 0$, we predict $y = 1$; if $\theta^T \mathbf{x} < 0$, then we predict $y = 0$. It determines the decision boundary



Decision boundary



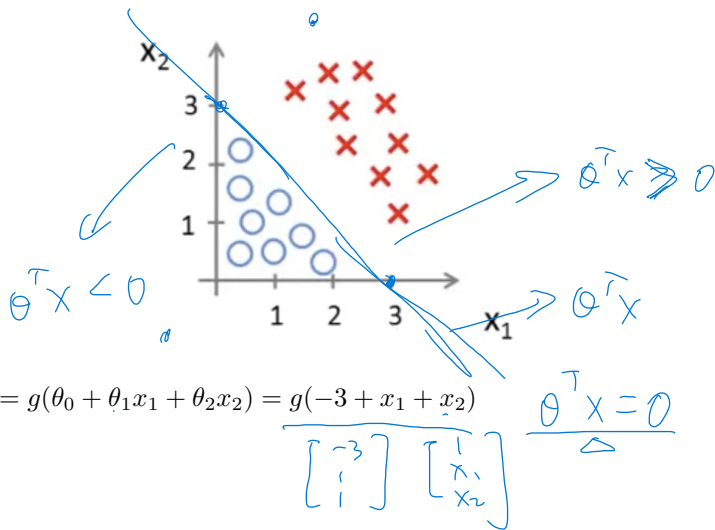
- In logistic regression, we have

$$h_{\theta}(\mathbf{x}) = g(\theta^{\top} \mathbf{x}) = P(y = 1 | \mathbf{x}; \theta) \in [0, 1], \quad (14)$$

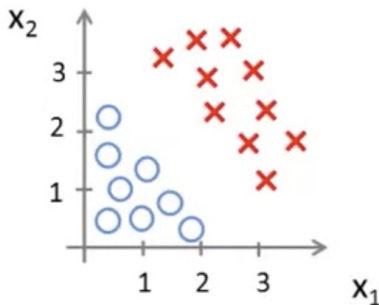
$$g(z) = \frac{1}{1 + \exp(-z)}, \quad (15)$$

- Suppose that if $h_{\theta}(\mathbf{x}) \geq 0.5$, then we predict $y = 1$; if $h_{\theta}(\mathbf{x}) < 0.5$, then we predict $y = 0$
- Correspondingly, if $\theta^{\top} \mathbf{x} \geq 0$, we predict $y = 1$; if $\theta^{\top} \mathbf{x} < 0$, then we predict $y = 0$. It determines the **decision boundary**
- **Decision boundary** is the curve/hyper-plane corresponding to $h_{\theta}(\mathbf{x}) = 0.5$, $\theta^{\top} \mathbf{x} = 0$

Decision boundary



Decision boundary



- $h_{\theta}(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = g(-3 + x_1 + x_2)$
- Predict $y = 1$ if $-3 + x_1 + x_2 \geq 0$ (plot above)

Decision boundary

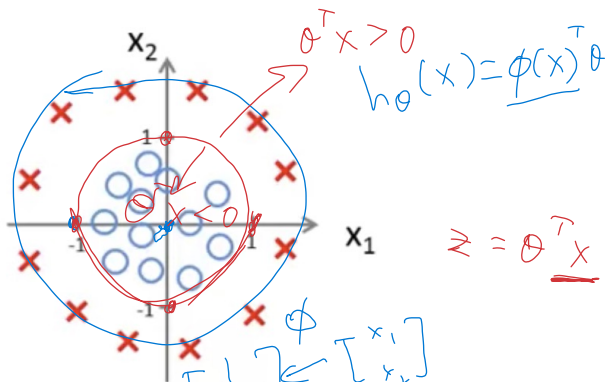


Figure: Non-linear decision boundary

$\theta^T \phi(x) =$
 $\bullet \ h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_1 x_1^2 + \theta_2 x_2^2) = g(\underline{-1 + x_1^2 + x_2^2})$
 $\theta^T x = 0$

Decision boundary

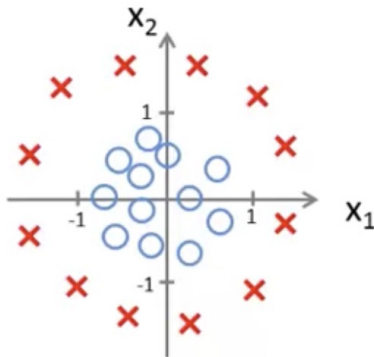


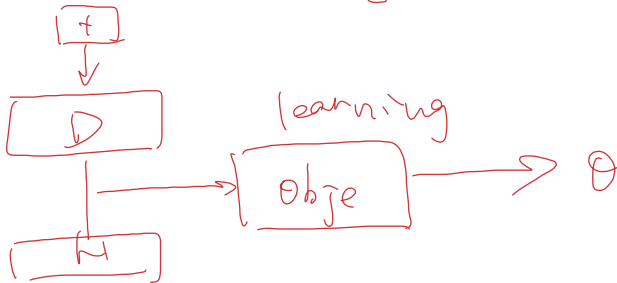
Figure: Non-linear decision boundary

- $h_{\theta}(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_1 x_1^2 + \theta_2 x_2^2) = g(-1 + x_1^2 + x_2^2)$
- Predict $y = 1$ if $-1 + x_1^2 + x_2^2 \geq 0$ (plot above)

- **Training set:** m training examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$

Cost function

- Training set: m training examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$
- Hypothesis function: $h_{\theta}(\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$



Cost function

- **Training set:** m training examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$
- **Hypothesis function:** $h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1+\exp(-\boldsymbol{\theta}^\top \mathbf{x})}$
- How to learn the model parameter $\boldsymbol{\theta}$?

Cost function

- **Training set:** m training examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$
- **Hypothesis function:** $h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1+\exp(-\boldsymbol{\theta}^\top \mathbf{x})}$
- How to learn the model parameter $\boldsymbol{\theta}$?
- We need to design a cost function/objective function

Cost function

Cost function

- Linear regression: $J(\theta) = \frac{1}{m} \sum_i^m \left(\underbrace{h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}} \right)^2$

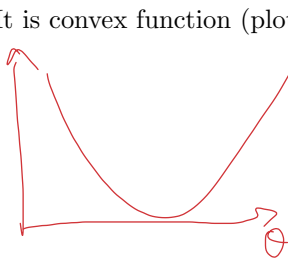
Cost function

- **Linear regression:** $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$
- $\text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}), y) = (h_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2$, which is called ℓ_2 **loss** or **residual sum of squares**

Cost function

Cost function

- **Linear regression:** $J(\theta) = \frac{1}{m} \sum_i^m \underbrace{(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2}_{= \theta^T \mathbf{x}}$
- $\text{cost}(h_{\theta}(\mathbf{x}), y) = (h_{\theta}(\mathbf{x}) - y)^2$, which is called ℓ_2 loss or residual sum of squares
- It is convex function (plot below) w.r.t. θ for linear regression



$$\begin{aligned} & (\theta^T \mathbf{x} - y)^2 \\ & 2(\theta^T \mathbf{x} - y) \cdot \mathbf{x} \\ & \underline{2 \mathbf{x}^T \mathbf{x} \geq 0} \end{aligned}$$

Cost function

$$\frac{2 h_{\theta}(x)}{2} = \frac{h \cdot (1-h)}{1 + e(-\theta^T x)}$$

Cost function

- **Linear regression:** $J(\theta) = \frac{1}{m} \sum_i^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2$
- $\text{cost}(h_{\theta}(\mathbf{x}), y) = (h_{\theta}(\mathbf{x}) - y)^2$, which is called ℓ_2 loss or residual sum of squares
- It is convex function (plot below) w.r.t. θ for linear regression
- However, it is non-convex (plot below) w.r.t. θ for logistic regression. Why? (please derive it)

$$\frac{\partial J(\theta)}{\partial \theta} = 2(h - y) \cdot h \cdot (1-h) \cdot (-x) = 2(h^2 - hy)(1-h)(-x)$$

$$\frac{\partial^2 J(\theta)}{\partial \theta^2} = 2[2h - 3h^2 - y - 2hy](-x) \cdot (1-h) \cdot h \cdot (-x) = 2(h^2 - h^3 - hy - h^2y)(-x)$$

Cost function

$$\underline{P} \rightarrow \underline{-\log p} \rightarrow \underline{\sum p (-\log p)}$$

- **Cross-entropy:**

$$H(p, q) = - \int_x \underline{p(x) \log(q(x))} d(x) \text{ or } \left[- \sum_x \underline{p(x) \log(q(x))} \right]$$

- We set

$$\underline{y(\mathbf{x}) = P(y = 1|\mathbf{x})}, \quad \underline{h_{\theta}(\mathbf{x}) = P(y = 1|\mathbf{x}; \theta)}$$

Cost function

- Cross-entropy:

$$H(p, q) \neq H(q, p)$$

$$H(p, q) = - \int_x p(x) \log(q(x)) d(x) \text{ or } - \sum_x p(x) \log(q(x))$$

- We set

$$\underline{y(\mathbf{x})} = \underbrace{P(y=1|\mathbf{x})}_{=\{0,1\}}, \underline{h_{\theta}(\mathbf{x})} = P(y=1|\mathbf{x}; \theta)$$

- Cross-entropy loss:

$$\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) = H(\underline{y(\mathbf{x})}, \underline{h_{\theta}(\mathbf{x})}) \quad (16)$$

$$= - \sum_{\mathbf{x}} y(\mathbf{x}) \log(h_{\theta}(\mathbf{x})) \quad (17)$$

$$= \begin{cases} -\log(h_{\theta}(\mathbf{x})), & \text{if } \underline{y(\mathbf{x})} = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})), & \text{if } \underline{y(\mathbf{x})} = 0 \end{cases} \quad (18)$$

Cost function

Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

- For $y = 1$, if $h_{\boldsymbol{\theta}}(\mathbf{x}) = 1$, i.e., $P(y = 1|\mathbf{x}; \boldsymbol{\theta}) = 1$, then the prediction equals to the ground-truth label, the cost is 0

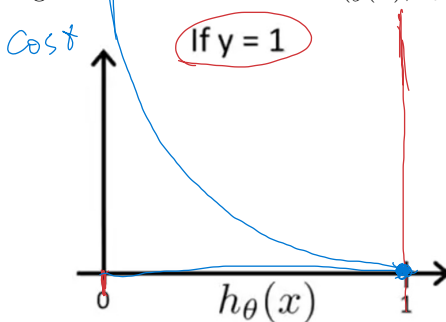
Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) = \begin{cases} -\log(h_{\theta}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

(Handwritten note: $\log(0)$ is written above the first case)

- For $y = 1$, if $h_{\theta}(\mathbf{x}) = 1$, i.e., $P(y = 1|\mathbf{x}; \theta) = 1$, then the prediction equals to the ground-truth label, the cost is 0
- For $y = 1$, if $h_{\theta}(\mathbf{x}) \rightarrow 0$, i.e., $P(y = 1|\mathbf{x}; \theta) \rightarrow 0$, then it should be penalized with a very large cost. Here we have $\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) \rightarrow \infty$



Cost function

Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

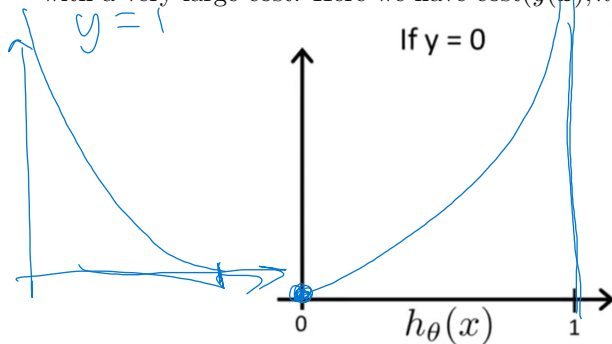
- For $y = 0$, if $h_{\boldsymbol{\theta}}(\mathbf{x}) = 0$, *i.e.*, $P(y = 1|\mathbf{x}; \boldsymbol{\theta}) = 0$, then the prediction equals to the ground-truth label, the cost is 0

Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

- For $y = 0$, if $h_{\boldsymbol{\theta}}(\mathbf{x}) = 0$, i.e., $P(y = 1|\mathbf{x}; \boldsymbol{\theta}) = 0$, then the prediction equals to the ground-truth label, the cost is 0
- For $y = 0$, if $h_{\boldsymbol{\theta}}(\mathbf{x}) \rightarrow 1$, i.e., $P(y = 1|\mathbf{x}; \boldsymbol{\theta}) \rightarrow 0$, then it should be penalized with a very large cost. Here we have $\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) \rightarrow \infty$



Cost function

- **Cross-entropy loss:**

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

Cost function

- Cross-entropy loss:

$$\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) = \begin{cases} -\log(h_{\theta}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

- Which states are true?

- If $h_{\theta}(\mathbf{x}) = y$, then $\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) = 0$ for both $y = 0$ and $y = 1$ ✓
- If $y = 0$, then $\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) \rightarrow \infty$ as $h_{\theta}(\mathbf{x}) \rightarrow 1$ ✓
- If $y = 0$, then $\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) \rightarrow \infty$ as $h_{\theta}(\mathbf{x}) \rightarrow 0$ ✗
- Regardless whether $y = 0$ or $y = 1$, if $h_{\theta}(\mathbf{x}) = 0.5$, then $\text{cost}(y(\mathbf{x}), h_{\theta}(\mathbf{x})) > 0$ ✓

Cost function of logistic regression

- Cost function of logistic regression

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m \text{cost}(y^{(i)}, h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})), \quad (19)$$

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases} \quad (20)$$

Cost function of logistic regression

- Cost function of logistic regression

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m \text{cost}(y^{(i)}, h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})), \quad (19)$$

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases} \quad (20)$$

- The above cost function can be simplified as follows

$$y^{(i)} \in \{0, 1\}$$

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_i^m [\underbrace{y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))}_{\text{if } y^{(i)}=1} + \underbrace{(1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))}_{\text{if } y^{(i)}=0}]. \quad (21)$$

Cost function of logistic regression

- Cost function of logistic regression

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m \text{cost}(y^{(i)}, h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})), \quad (19)$$

$$\text{cost}(y(\mathbf{x}), h_{\boldsymbol{\theta}}(\mathbf{x})) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})), & \text{if } y(\mathbf{x}) = 0 \end{cases} \quad (20)$$

- The above cost function can be simplified as follows

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))]. \quad (21)$$

Handwritten note: $\frac{h(z)}{\Delta z} = \frac{h \cdot (1-h)}{\Delta}$

- It is convex function (plot below) w.r.t. $\boldsymbol{\theta}$. (Please derive it)

Gradient descent

$$J(\theta) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]. \quad (22)$$

- Learning θ by $\min_{\theta} J(\theta)$
- **Gradient descent:** repeat the following update until convergence

$$\theta \leftarrow \theta - \eta \nabla_{\theta} J(\theta) \quad (23)$$

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_i^m [h_{\theta}(x^{(i)}) - y^{(i)}] x^{(i)} \quad (24)$$

- How to define convergence? Calculating the changes of $J(\theta)$ or θ in the last K steps, if the change is lower than a threshold, then it can be seen as convergence. Remember that choosing suitable learning rate η is important to achieve a good converged solution.

10^{-4}

Gradient descent

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))]. \quad (22)$$

- Learning $\boldsymbol{\theta}$ by $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- **Gradient descent:** repeat the following update until convergence

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (23)$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m [h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}^{(i)} \quad (24)$$

Gradient descent

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))]. \quad (22)$$

- Learning $\boldsymbol{\theta}$ by $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- **Gradient descent:** repeat the following update until convergence

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (23)$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m [h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}^{(i)} \quad (24)$$

- How to define convergence?

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))]. \quad (22)$$

- Learning $\boldsymbol{\theta}$ by $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- **Gradient descent:** repeat the following update until convergence

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (23)$$

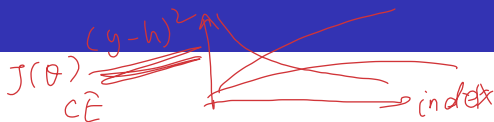
$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i^m [h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}^{(i)} \quad (24)$$

- How to define convergence? Calculating the changes of $J(\boldsymbol{\theta})$ or $\boldsymbol{\theta}$ in the last K steps, if the change is lower than a threshold, then it can be seen as convergence. Remember that choosing suitable learning rate η is important to achieve a good converged solution.

Gradient descent

Suppose you are running a logistic regression model, and you should observe the learning procedure to find a suitable learning rate η . Which of the following is reasonable to make sure η is set properly and that the gradient descent is running correctly?

Gradient descent



Suppose you are running a logistic regression model, and you should observe the learning procedure to find a suitable learning rate η . Which of the following is reasonable to make sure η is set properly and that the gradient descent is running correctly?

- Plot $J(\theta) = -\frac{1}{m} \sum_i^m (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))^2$ as a function of the number of iterations (*i.e.*, the horizontal axis is the iteration number) and make sure $J(\theta)$ is decreasing on every iteration. ✗
- Plot $J(\theta) = -\frac{1}{m} \sum_i^m [y^{(i)} \log(h_{\theta}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))]$ as a function of the number of iterations (*i.e.*, the horizontal axis is the iteration number) and make sure $J(\theta)$ is decreasing on every iteration. ✓
- Plot $J(\theta)$ as a function of θ and make sure it is decreasing on every iteration.
- Plot $J(\theta)$ as a function of θ and make sure it is convex.

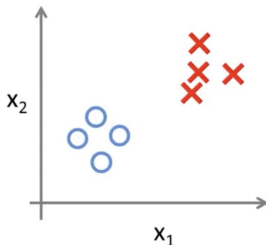
Multi-class classification

- In above examples and derivations, we only consider the binary classification problem, *i.e.*, $y \in \{0, 1\}$.

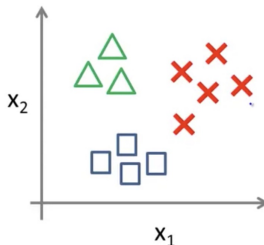
Multi-class classification

- In above examples and derivations, we only consider the binary classification problem, *i.e.*, $y \in \{0, 1\}$.
- However, many practical problems involve with multi-class classification:
 - Whether forecast: sunny, cloudy, rain, snow
 - Email tagging: work, friends, families, hobby

Binary classification:

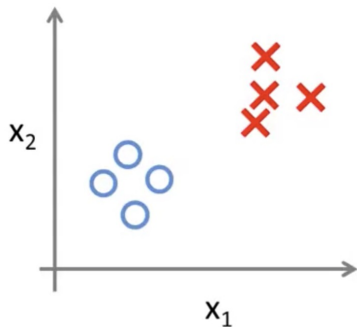


Multi-class classification:

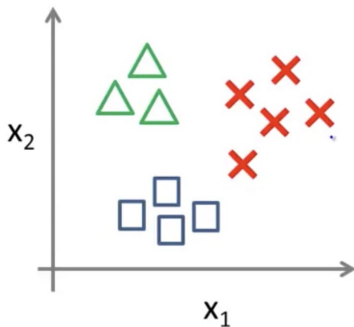


Multi-class classification

Binary classification:

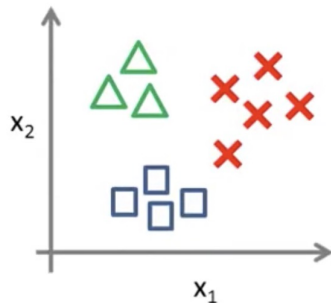



Multi-class classification:





Multi-class classification: one-vs-all

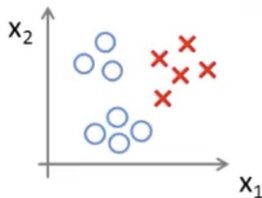
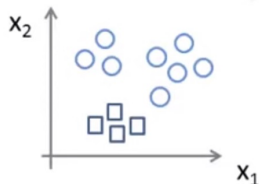
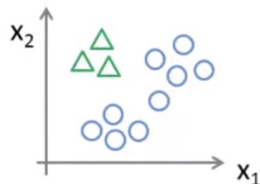
One-vs-all (one-vs-rest):



Class 1: 

Class 2: 

Class 3: 



Multi-class classification: one-vs-all

One-vs-all logistic regression:

- Train a binary logistic regression $h_{\theta_i}(\cdot)$ for each class i , by setting all samples of other classes as negative class
- For a new testing sample \mathbf{x} , predict its class as $\arg \max_i h_{\theta_i}(\mathbf{x})$.

Multi-class classification: one-vs-all

One-vs-all logistic regression:

- Train a binary logistic regression $h_{\theta_i}(\cdot)$ for each class i , by setting all samples of other classes as negative class
- For a new testing sample \mathbf{x} , predict its class as $\arg \max_i h_{\theta_i}(\mathbf{x})$.

Pros: Easy to implement

Multi-class classification: one-vs-all

One-vs-all logistic regression:

- Train a binary logistic regression $h_{\theta_i}(\cdot)$ for each class i , by setting all samples of other classes as negative class
- For a new testing sample \mathbf{x} , predict its class as $\arg \max_i h_{\theta_i}(\mathbf{x})$.

Pros: Easy to implement

Cons: The training cost is too high, and is difficult to scale to tasks with large number of classes.

Multi-class classification: logistic regression with softmax function

- Softmax function:

$$h_{\boldsymbol{\theta}_i}(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_i^\top \mathbf{x})}{\sum_i \exp(\boldsymbol{\theta}_i^\top \mathbf{x})} = P(y = i | \mathbf{x}; \boldsymbol{\Theta}), \quad (25)$$

where $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_i\}_i^C$ with C being the number of classes.

Multi-class classification: logistic regression with softmax function

- Softmax function:

$$h_{\boldsymbol{\theta}_i}(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_i^\top \mathbf{x})}{\sum_i \exp(\boldsymbol{\theta}_i^\top \mathbf{x})} = P(y = i | \mathbf{x}; \boldsymbol{\Theta}), \quad (25)$$

where $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_i\}_i^C$ with C being the number of classes.

- Cost function:

$$J(\boldsymbol{\Theta}) = -\frac{1}{m} \sum_j^m \sum_i^C [\mathbb{I}(\underline{y^{(j)} = i}) \log(h_{\boldsymbol{\theta}_i}(\mathbf{x}^{(j)}))], \quad (26)$$

where $\mathbb{I}(\underline{a}) = 1$ if a is true, otherwise $\mathbb{I}(a) = 0$.

Multi-class classification: logistic regression with softmax function

- Softmax function:

$$h_{\boldsymbol{\theta}_i}(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_i^\top \mathbf{x})}{\sum_i \exp(\boldsymbol{\theta}_i^\top \mathbf{x})} = P(y = i | \mathbf{x}; \boldsymbol{\Theta}), \quad (25)$$

where $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_i\}_i^C$ with C being the number of classes.

- Cost function:

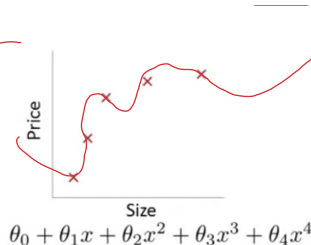
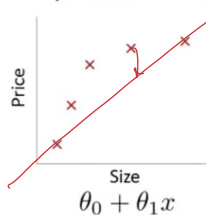
$$J(\boldsymbol{\Theta}) = -\frac{1}{m} \sum_j^m \sum_i^C [\mathbb{I}(y^{(j)} = i) \log(h_{\boldsymbol{\theta}_i}(\mathbf{x}^{(j)}))], \quad (26)$$

where $\mathbb{I}(a) = 1$ if a is true, otherwise $\mathbb{I}(a) = 0$.

- It can also be optimized by gradient descent. (Please derive its gradient)

Overfitting in linear regression

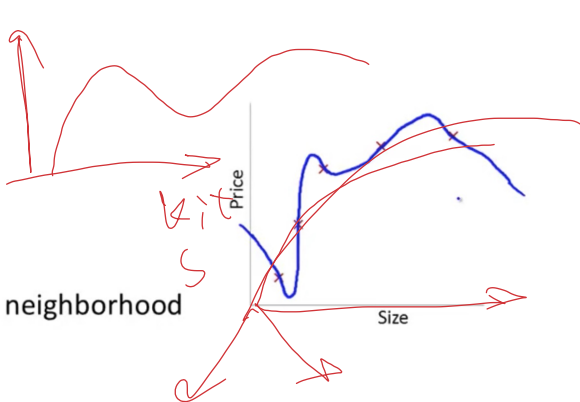
Example: Linear regression (housing prices)



Overfitting in linear regression

Addressing overfitting:

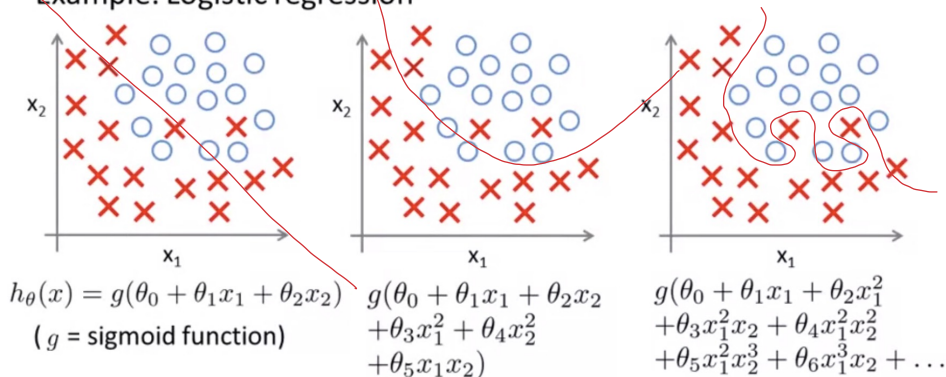
x_1 = size of house
 x_2 = no. of bedrooms
 x_3 = no. of floors
 x_4 = age of house
 x_5 = average income in neighborhood
 x_6 = kitchen size
 \vdots
 x_{100}



Overfitting: If we have too many features, the learned hypothesis may fit the training data very well (low bias), but fail to generalize to new examples.

Overfitting in logistic regression

Example: Logistic regression



Addressing Overfitting

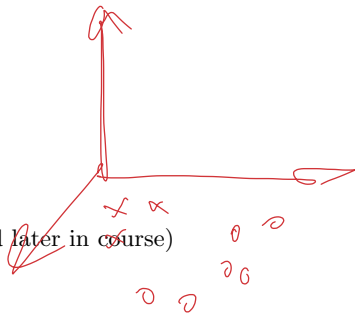
Two options:

Addressing Overfitting

Two options:

① Reducing the number of features:

- Feature selection
- Dimensionality reduction (introduced later in course)



Addressing Overfitting

Two options:

① Reducing the number of features:

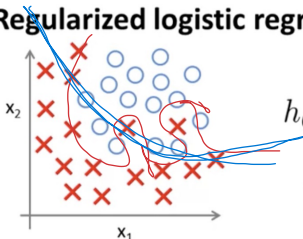
- Feature selection
- Dimensionality reduction (introduced later in course)

② Regularization:

- Keep all features, but reduce magnitude/value of each parameter, such that each feature contributes a bit to predict y

Regularized logistic regression

Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

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

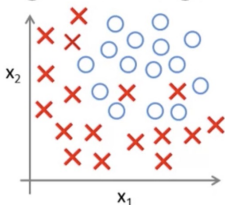
$$\frac{+2\lambda}{m} \|\theta\|_2^2$$

↓

$$\sum_i \theta_i^2$$

Regularized logistic regression

Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

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

$$J(\theta) = - \frac{1}{m} \sum_i \left[y^{(i)} \log(h_{\theta}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2.$$

Regularized logistic regression

$$J(\theta) = CE(\theta)$$

Gradient descent

Repeat {

$$\underbrace{\theta_j}_{\theta} := \underbrace{\theta_j} - \underbrace{\alpha}_{\uparrow} \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{(j=0,1,2,3,\dots,n)}$$

θ_0

Regularized logistic regression

Gradient descent

Repeat {

$$\begin{aligned} \rightarrow \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \rightarrow \theta_j &:= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &\quad (j = \underline{1, 2, 3, \dots, n}) \end{aligned}$$

Note: the bias parameter θ_0 is not penalized.

$$\begin{aligned} J(\theta) &= CE(\theta) \\ &\quad + \frac{\lambda}{2} \|\theta\|^2 \end{aligned}$$

$\sum_{j=1}^n \theta_j^2$

Regularized logistic regression

$$J(\theta) = \underbrace{CE(\theta)}_{\text{index}} + \frac{\lambda}{2} \|\theta\|^2$$

When using regularized logistic regression, which of these is the best way to monitor whether gradient descent is working correctly?

- ☐ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right]$ as a function of the number of iterations and make sure it's decreasing. X
- ☐ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] - \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ as a function of the number of iterations and make sure it's decreasing. X
- ☒ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ as a function of the number of iterations and make sure it's decreasing. ✓
- ☐ Plot $\sum_{j=1}^n \theta_j^2$ as a function of the number of iterations and make sure it's decreasing. X

Generalized linear regression

- Linear model:

$$\left\{ \begin{array}{l} \mu(\mathbf{x}|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \phi(\mathbf{x}), \\ y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \end{array} \right. \quad \begin{array}{l} (27) \\ (28) \end{array}$$

where f denotes a distribution function.

Generalized linear regression

- **Linear model:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \phi(\mathbf{x}), \quad (27)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (28)$$

where f denotes a distribution function.

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad (29)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (30)$$

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

Generalized linear regression

- **Linear model:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}), \quad (27)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (28)$$

where f denotes a distribution function.

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

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = g^{-1}(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})), \quad (29)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim f(\mu(\mathbf{x}|\boldsymbol{\theta})), \quad (30)$$

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 $g(a) = a$.

Logistic regression: probabilistic modeling

- Logistic regression:

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

$$\ln \frac{u}{1-u} \quad (31)$$

$$\begin{cases} \mu(\mathbf{x}|\boldsymbol{\theta}) = \text{Sigmoid}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \\ y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \end{cases}$$

(32)

$$\arg \max_{\boldsymbol{\theta}} \left[\log \left[u^y \cdot (1-u)^{(1-y)} \right] \right]$$

$$\min_{\boldsymbol{\theta}} C \bar{E}(\boldsymbol{\theta}) = y \log u + (1-y) \log(1-u)$$

$$- C \bar{E}(\boldsymbol{\theta}) = y \log \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}}} + (1-y) \log \frac{e^{-\boldsymbol{\theta}^\top \mathbf{x}}}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}}}$$

Logistic regression: probabilistic modeling

- **Logistic regression:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \text{Sigmoid}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad (31)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad (32)$$

- We have

$$P(y|\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} \mu & \text{if } y = 1 \\ 1 - \mu & \text{if } y = 0 \end{cases} \quad (33)$$

Logistic regression: probabilistic modeling

- **Logistic regression:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \text{Sigmoid}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad (31)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad (32)$$

- We have

$$P(y|\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} \mu & \text{if } y = 1 \\ 1 - \mu & \text{if } y = 0 \end{cases} \quad (33)$$

- The log-likelihood is

$$\mathcal{L}(\boldsymbol{\theta}) = y \log(\mu) + (1 - y) \log(1 - \mu) \quad (34)$$

Logistic regression: probabilistic modeling

- **Logistic regression:**

$$\mu(\mathbf{x}|\boldsymbol{\theta}) = \text{Sigmoid}(\boldsymbol{\theta}^\top \phi(\mathbf{x})), \quad (31)$$

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad (32)$$

- We have

$$P(y|\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} \mu & \text{if } y = 1 \\ 1 - \mu & \text{if } y = 0 \end{cases} \quad (33)$$

- The log-likelihood is

$$\mathcal{L}(\boldsymbol{\theta}) = y \log(\mu) + (1 - y) \log(1 - \mu) \quad (34)$$

- Thus, we obtain

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \equiv \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (35)$$

Logistic regression: probabilistic modeling

- **Logistic regression:**

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

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad (37)$$

Logistic regression: probabilistic modeling

- Logistic regression:

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

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad \text{Lap}(\boldsymbol{\theta}|b) \quad (37)$$

- Regularized Logistic regression:** we assume $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \sigma^2\mathbf{I})$, then we have

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \log \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \sigma^2\mathbf{I}) \equiv \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}_{\sigma^2} \quad (38)$$

$$\frac{\sum |\theta_j|}{\sigma} \exp\left(-\frac{|x-a|}{b}\right)$$

Logistic regression: probabilistic modeling

- **Logistic regression:**

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

$$y(\mathbf{x}|\boldsymbol{\theta}) \sim \text{Bernoulli}(\mu(\mathbf{x}|\boldsymbol{\theta})). \quad (37)$$

- **Regularized Logistic regression:** we assume $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \sigma^2\mathbf{I})$, then we have

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \log \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \sigma^2\mathbf{I}) \equiv \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad (38)$$

- **Regularized Logistic regression:** if we assume $\boldsymbol{\theta} \sim \text{Laplace}(\boldsymbol{\theta}|\mathbf{0}, b)$, then we have

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \log \text{Laplace}(\boldsymbol{\theta}|\mathbf{0}, b) \equiv \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n |\theta_j|} \quad (39)$$

Summary of linear models

- Linear model is the linear function of the parameter θ , rather than the input feature

Summary of linear models

- 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

Summary of linear models

- 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
- Logistic model is also a special case of generalized linear model. It is used to solve classification task because the output range of its hypothesis function (*i.e.*, sigmoid function) is $[0, 1]$, which can be seen as posterior probability

Summary of linear models

$$\theta^T \phi(x)$$

- 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
- Logistic model is also a special case of generalized linear model. It is used to solve classification task because the output range of its hypothesis function (*i.e.*, sigmoid function) is $[0, 1]$, which can be seen as posterior probability
- Different variants of linear models correspond to different distributions of $p(y|x, \theta)$ and $p(\theta)$, according to the task and the data, *i.e.*, handling outliers or alleviating overfitting

$p(y x, \theta)$	$p(\theta)$
Gau	unit ()
B	Gau \rightarrow
La	La \rightarrow