

LINEAR MODEL

Bùi Tiến Lên

2023



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

Contents



1. Linear Regression
2. Classification
3. Logistic Regression
4. Softmax Regression
5. Capacity, Overfitting and Underfitting



Notation

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

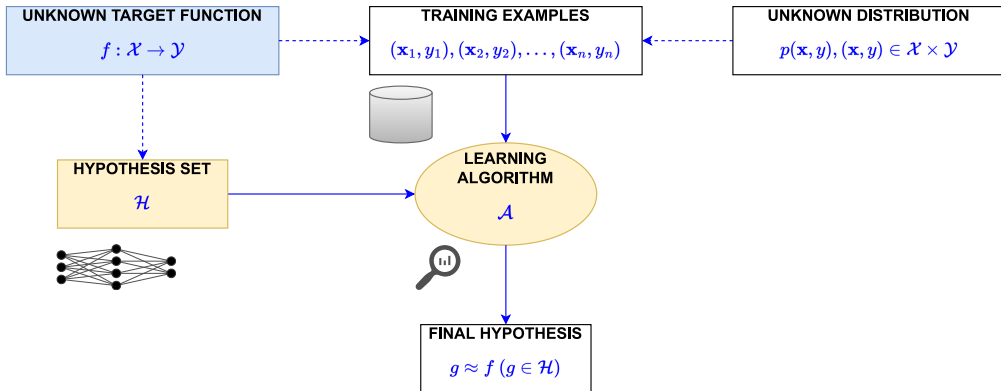
Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

symbol	meaning
$a, b, c, N \dots$	scalar number
$\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$	column vector
$\mathbf{X}, \mathbf{Y} \dots$	matrix
\mathbb{R}	set of real numbers
\mathbb{Z}	set of integer numbers
\mathbb{N}	set of natural numbers
\mathbb{R}^D	set of vectors
$\mathcal{X}, \mathcal{Y}, \dots$	set
\mathcal{A}	algorithm

operator	meaning
\mathbf{w}^T	transpose
$\mathbf{X}\mathbf{Y}$	matrix multiplication
\mathbf{X}^{-1}	inverse



Learning diagram



Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization



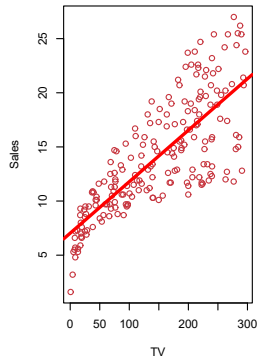
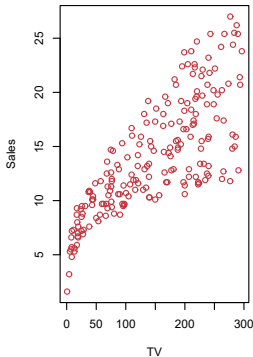
Linear Regression

- Simple Linear Model
- Weighted Linear Model
- Linear Basis Function Model

Problem 1



- Consider the Advertising data set \mathcal{D}_{train} consists of the **sales** of that product in 200 different markets, along with advertising budgets for the product in each of those markets for the media **TV**. Find the *relationship* between **TV** (input) and **sales** (output)



Linear Regression Model



Concept 1

A **linear regression** is a model that assumes a **linear relationship** between **inputs** and the **output**.



Problem Statement

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- **The requirement** is to build a system that can take a vector $\mathbf{x} \in \mathbb{R}^{D+1}$ as **input** and **predict** the value of a scalar $y \in \mathbb{R}$ as its **output**
- **The hypothesis set** \mathcal{H}

$$y \approx \hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \quad (1)$$

where \hat{y} be the value that our model (function) predicts y and $\mathbf{w} \in \mathbb{R}^{D+1}$ is a *vector of parameters* of the model



Problem Statement (cont.)

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- **Task T :** to predict y from \mathbf{x} by outputting $\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- The train set \mathcal{D}_{train} denoted as (\mathbf{X}, \mathbf{y}) including N samples $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N)\}$, construct the matrix \mathbf{X} and the vectors \mathbf{y} and $\hat{\mathbf{y}}$

$$\underbrace{\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}, \quad \underbrace{\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}}_{\text{output vector}} \quad (2)$$



Problem Statement (cont.)

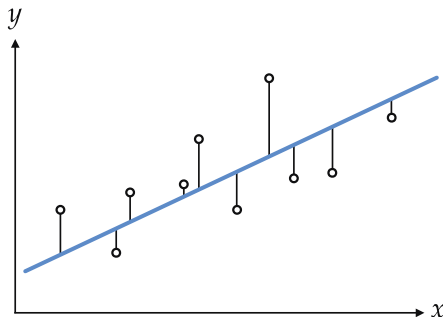
- Performance measure P :

Concept 2

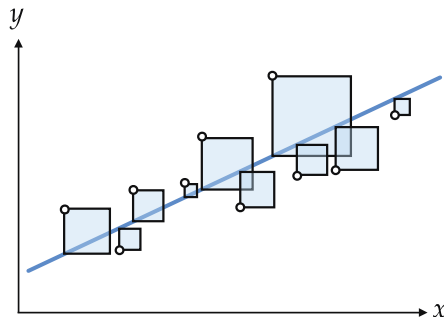
The *mean squared error* MSE_{train} of the model on the train set \mathcal{D}_{train}

$$MSE_{train} = \frac{1}{N} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2 \quad (3)$$

Problem Statement (cont.)



MAE



MSE



Problem Statement (cont.)

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- **The learning goal:** find the vector of parameter \mathbf{w} such that

$$\mathbf{w} = \arg \min_{\mathbf{w}} (MSE_{train}) \quad (4)$$



Solving Problem

Solution

- Compute the gradient of MSE_{train}

$$\begin{aligned}\nabla_{\mathbf{w}}(MSE_{train}) &= \nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) \\ &= 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y}\end{aligned}\quad (5)$$

- If MSE_{train} reach the min value then $\nabla_{\mathbf{w}}(MSE_{train}) = 0$

$$\begin{aligned}\nabla_{\mathbf{w}}(MSE_{train}) &= 0 \\ \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} &= 0 \\ \mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}\quad (6)$$





Solving Problem (cont.)

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

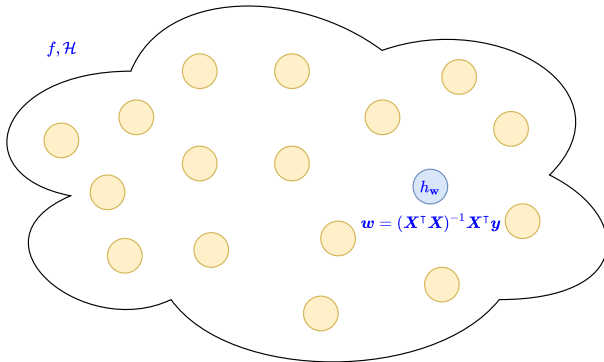
Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization



Programming Example

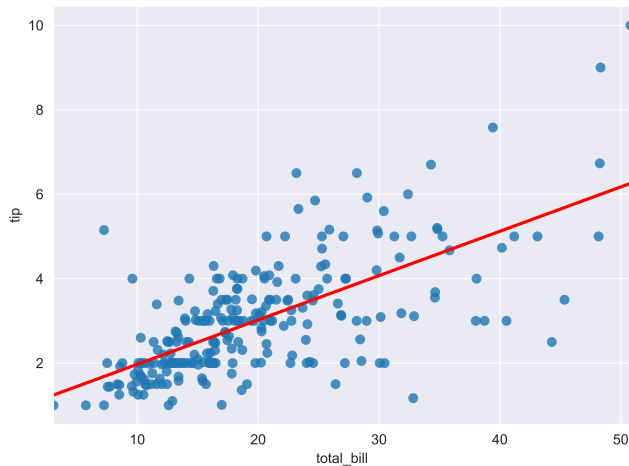


- Use seaborn to read tips dataset and find the linear relationship between total_bill and tip

```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

sns.set_style("darkgrid")
tips = sns.load_dataset("tips")
sns.regplot(x="total_bill", y="tip", data=tips, ci=None, line_kws=
           ={'color': 'red'})
plt.show()
```

Programming Example (cont.)



Word Example



1. Find the linear regression function $y = f(x) = w_0 + w_1x$ given the following data set \mathcal{D}

input x	target y
1	2
2	3
3	3
4	5

2. Find the linear regression function $y = f(\mathbf{x}) = f(x_1, x_2) = w_0 + w_1x_1 + w_2x_2$ given the following data set \mathcal{D}

input \mathbf{x}	target y
(1, 1)	1
(2, 3)	3
(3, 4)	4
(4, 3)	5



Discussion

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- D is a large number
- Online learning
- Limitations of **the model (hypothesis set)**



Weighted Linear Model

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- In some cases the observations may be weighted; for example, they may not be equally reliable. In this case, we find the vector of parameters \mathbf{w} to minimize the weighted sum of squares of errors

$$E_{train} = \sum_{n=1}^N a_n (\hat{y}_n - y_n)^2 \quad (7)$$



Solving Problem

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

1. Construct the matrices \mathbf{X} , \mathbf{A} and the vector \mathbf{y}

$$\underbrace{\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \underbrace{\mathbf{A} = \begin{bmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_N \end{bmatrix}}_{\text{weight matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}} \quad (8)$$

2. Calculate the vector of parameters

$$\mathbf{w} = (\mathbf{X}^T \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{A} \mathbf{y} \quad (9)$$



Linear in What?

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- **Linearity in the weights**

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_Dx_D \quad (10)$$



Concept 3

A **linear basis function model** is a linear combination of **fixed nonlinear functions** of the input variables

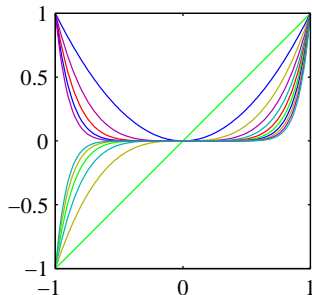
$$h_{\mathbf{w}}(\mathbf{x}) = w_0\phi_0(\mathbf{x}) + w_1\phi_1(\mathbf{x}) + \dots + w_M\phi_M(\mathbf{x}) \quad (11)$$

where $\phi_i(\mathbf{x})$ are basis functions

Some types of basis functions



- Polynomial

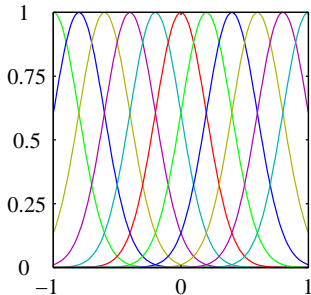


$$\phi_j(x) = x^j \quad (12)$$



Some types of basis functions (cont.)

- Gaussian

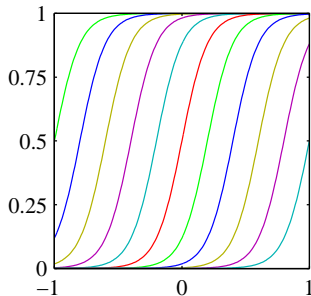


$$\phi_j(x; \mu_j, s_j) = \exp\left(-\frac{(x - \mu_j)^2}{s_j^2}\right) \quad (13)$$



Some types of basis functions (cont.)

- Sigmoid

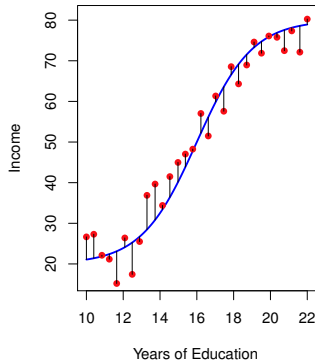
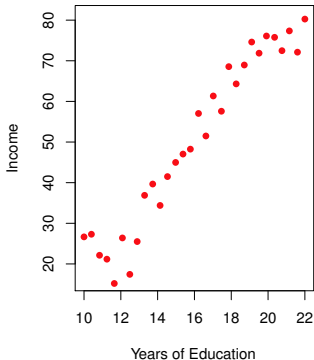


$$\phi_j(x; \mu_j, s_j) = \frac{1}{1 + e^{-\frac{x - \mu_j}{s_j}}} \quad (14)$$

Problem 2



- Find the relationship between Years of Education and Income based on the given data





Problem Statement

- **The requirement** is to build a system that can take a vector $\mathbf{x} \in \mathbb{R}^D$ as **input** and **predict** the value of a scalar $y \in \mathbb{R}$ as its **output**
- The hypothesis set \mathcal{H}

$$y \approx \hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) \quad (15)$$

where \hat{y} be the value that our model (function) predicts y , $\mathbf{w} \in \mathbb{R}^{M+1}$ is a *vector of parameters* of the model and $\boldsymbol{\phi}$ is a set of $M + 1$ basis functions

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \phi_0(\mathbf{x}) \\ \phi_1(\mathbf{x}) \\ \vdots \\ \phi_M(\mathbf{x}) \end{bmatrix} \quad (16)$$



Problem Statement (cont.)

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- **Task T :** to predict y from \mathbf{x} by outputting $\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$

- **Performance measure P :**

The *mean squared error* MSE_{train} of the model on the train set \mathcal{D}_{train} including N samples $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N)\}$

- **The learning goal:** find the vector of parameter \mathbf{w} such that

$$\mathbf{w} = \arg \min_{\mathbf{w}} (MSE_{train})$$



Solving Problem

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

1. Construct the matrix Φ and the vector \mathbf{y}

$$\underbrace{\Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{bmatrix}}_{\text{design matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}} \quad (17)$$

2. Calculate the vector of parameters

$$\mathbf{w} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} \quad (18)$$

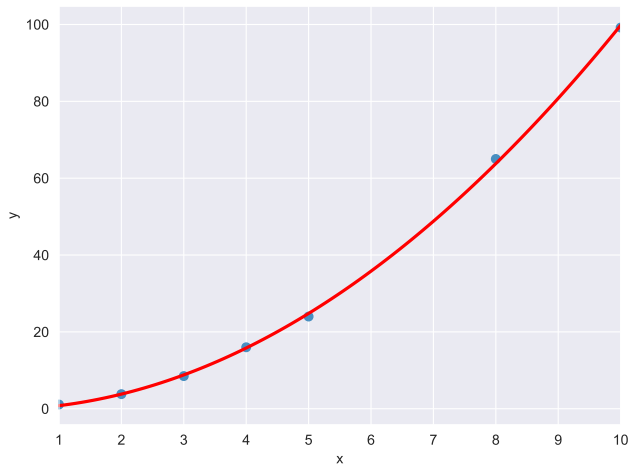
Programming Example



```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
sns.set_style("darkgrid")
x = [1, 2, 3, 4, 5, 8, 10]
y = [1.1, 3.8, 8.5, 16, 24, 65, 99.2]
sns.regplot(x, y, order=2, ci=None, line_kws={'color':'red'})
plt.xlabel('x')
plt.ylabel('y')
plt.show()
```

Programming Example (cont.)





Word Example

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

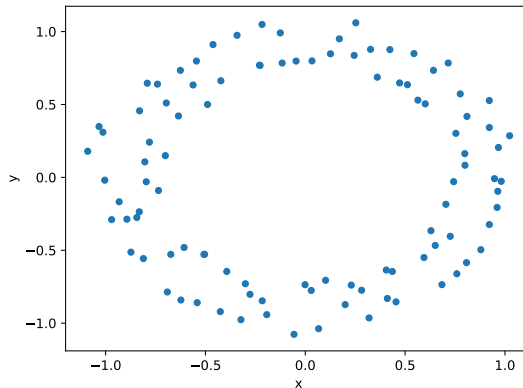
- Find a polynomial regression function $y = f(x) = w_0 + w_1x + w_2x^2$ given the data set \mathcal{D}

input x	target y
1	2
2	3
3	3
4	5

Puzzle



- What basis functions?





Classification



A real data set

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

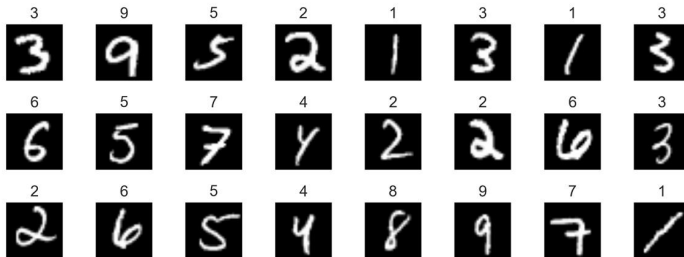
Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Some 16-by-16 pixel grayscale image from the MNIST database



Input representation



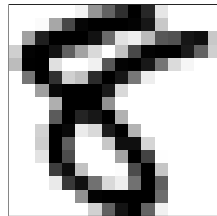
Input representation or feature extraction

- “raw” input

$$\begin{array}{ll} \text{pixels} & \mathbf{x}^T = (x_0 \quad x_1 \quad \dots \quad x_{256}) \\ \text{linear model} & \mathbf{w}^T = (\mathbf{w}_0 \quad \mathbf{w}_1 \quad \dots \quad \mathbf{w}_{256}) \end{array}$$

- **Feature extraction:** extract useful information

$$\begin{array}{ll} \text{intensity and symmetry} & \mathbf{x}^T = (x_1 \quad x_2) \\ \text{linear model} & \mathbf{w}^T = (\mathbf{w}_1 \quad \mathbf{w}_2) \end{array}$$



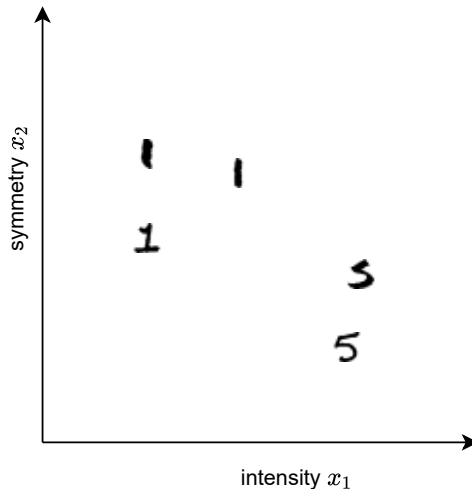
Classification

Logistic Regression

Softmax Regression

Capacity, Overfitting and Underfitting

Illustration of Features



The Problem with Categorical Data



- Some algorithms can work with categorical data directly.
 - For example, a decision tree can be learned directly from categorical data with no data transform required
- Many machine learning algorithms cannot operate on label data directly.
 - They require all input variables and output variables to be numeric.
- There are two common types of conversion: integer encoding and one-hot encoding



Integer Encoding

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- For categorical variables where no such ordinal relationship exists, the integer encoding is not enough.

id	color
1	red
2	green
3	blue
4	red
...	...

Integer encoding

id	color
1	1
2	2
3	3
4	1
...	...

One-Hot Encoding



- One-hot encoding ensures that machine learning does not assume that higher numbers are more important.

id	color
1	red
2	green
3	blue
4	red
...	...

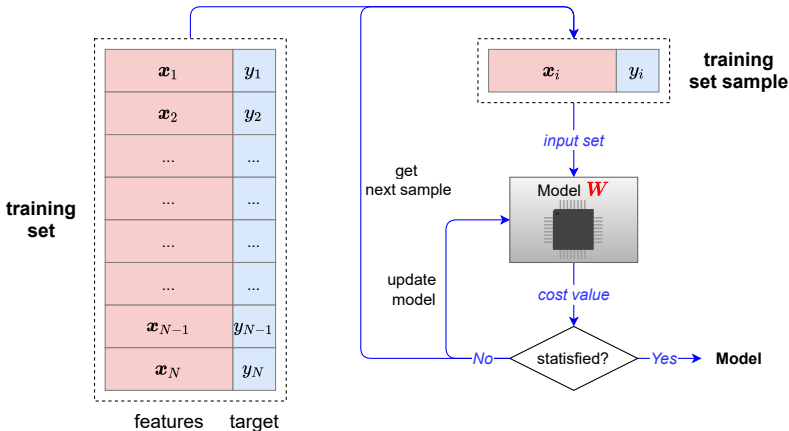
One-hot
encoding

id	color_red	color_green	color_blue
1	1	0	0
2	0	1	0
3	0	0	1
4	1	0	0
...

Classifier Training



- **Select** the learning model for **classifier**, e.g., Perceptron
- **Train** the classifier/model using a training set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$





Logistic Regression

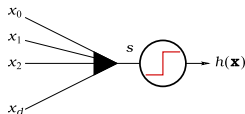
- Binary Classification
- Evaluation
- Multi-class Classification

A Third Linear Model

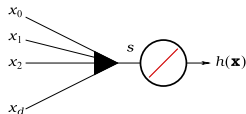


$$s = \sum_{i=0}^d w_i x_i$$

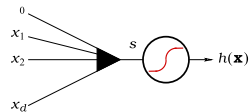
linear classification
 $h(\mathbf{x}) = \text{sign}(s)$



linear regression
 $h(\mathbf{x}) = s$



logistic regression
 $h(\mathbf{x}) = \sigma(s)$





The logistic function

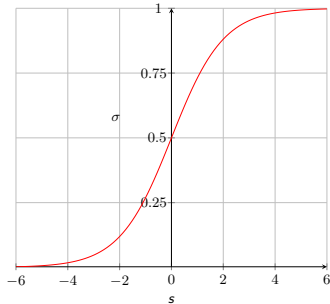
- The formula

$$\sigma(s) = \frac{1}{1 + e^{-s}} \quad (19)$$

- The **logistic function** converts a score to a probability
- Properties

$$\sigma(-s) = 1 - \sigma(s)$$

$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$





Probability Interpretation

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- $h(\mathbf{x}) = \sigma(s)$ can be interpreted as a probability
- For example, prediction of heart attacks
 - Input \mathbf{x} : cholesterol level, age, weight, etc.
 - The signal $s = \mathbf{w}^T \mathbf{x}$: risk score
 - $\sigma(s)$: probability of a heart attack



score \longrightarrow probability of heart attack



Problem Statement

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- The target function f is the probability distribution

$$f : \mathbb{R}^D \rightarrow [0, 1]$$

- Hypothesis set $h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$ and the conditional probability

$$P(y \mid \mathbf{x}, \mathbf{w}) = \begin{cases} h_{\mathbf{w}}(\mathbf{x}) & \text{for } y = 1 \\ 1 - h_{\mathbf{w}}(\mathbf{x}) & \text{for } y = 0 \end{cases} \quad (20)$$



Error measure

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

We define error measurement based on likelihood

- For each (\mathbf{x}, y) , y is generated by probability $h_{\mathbf{w}}(\mathbf{x})$.
- Plausible error measure based on **likelihood** of y given \mathbf{x} and \mathbf{w}

$$P(y \mid \mathbf{x}, \mathbf{w}) = z^y (1 - z)^{1-y} \quad (21)$$

where

$$z = h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) \quad (22)$$

- Likelihood of $\mathcal{D} = \{(\mathbf{x}_1, y_1) \dots (\mathbf{x}_N, y_N)\}$ given \mathbf{w} is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n, \mathbf{w}) \quad (23)$$



Error measure (cont.)

- **Learning goal:** maximizing likelihood

$$\begin{aligned} & \text{Maximize} && \prod_{n=1}^N P(y_n | \mathbf{x}_n, \mathbf{w}) \\ \Leftrightarrow & \text{Minimize} && -\log \prod_{n=1}^N P(y_n | \mathbf{x}_n, \mathbf{w}) \\ \Leftrightarrow & \text{Minimize} && -\sum_{n=1}^N (y_n \log z_n + (1 - y_n) \log(1 - z_n)) \end{aligned} \quad (24)$$

1 +



Learning Algorithm (Gradient Descent)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

1. Initialize the weights (parameters) at $t = 0$ \mathbf{w}_0

2. For $t = 1, 2, 3, \dots$ do

2.1 Compute the outputs z_n for each \mathbf{x}_n ($n = 1, \dots, N$)

$$z_n = \sigma(\mathbf{w}_t^T \mathbf{x}_n) \quad (25)$$

2.2 Compute the gradient

$$\nabla_{\mathbf{w}} E = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n (z_n - y_n) \quad (26)$$

2.3 Update the weights

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} E \quad (27)$$

where η is a learning rate (*hyper-parameter*)

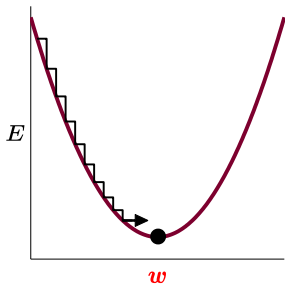
Iterate the next step until \mathbf{w} is not changes

3. Return the final weights \mathbf{w}

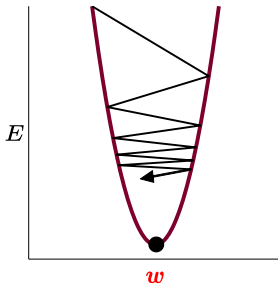


Learning Rate

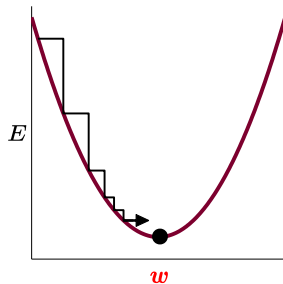
- How η affects the algorithm?



η too small



η too large



η right

Programming Example



```
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="darkgrid")

# Load the example titanic dataset
df = sns.load_dataset("titanic")

# Make a custom palette with gendered colors
pal = dict(male="#6495ED", female="#F08080")

# Show the survival probability as a function of age and sex
g = sns.lmplot(x="age", y="survived", col="sex", hue="sex", data=df,
               palette=pal, y_jitter=.02, logistic=True, ci=None)
g.set(xlim=(0, 80), ylim=(-.05, 1.05))
plt.show()
```



Programming Example (cont.)

Linear Regression

- Simple Linear Model
- Weighted Linear Model
- Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

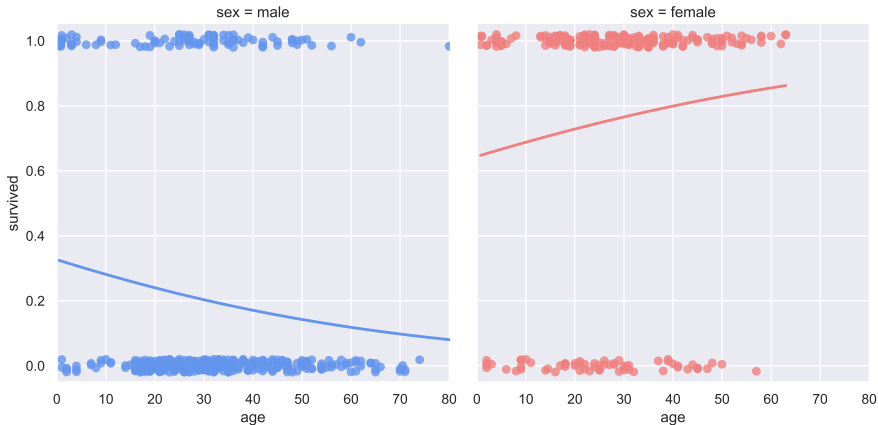
- Evaluation
- Multi-class Classification

Softmax Regression

- Softmax Regression
- Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

- Model Capacity
- Model vs. Data
- Bias-Variance
- Tradeoff of Capacity
- Regularization
- Tradeoff of Regularization





Evaluation

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Consider two-class problem with two classes \oplus and \ominus
- The performance of logistic regression model is based on a threshold th

$$\begin{aligned} y \text{ is } \oplus & \text{ if } P(y \mid \mathbf{x}) \geq th \\ y \text{ is } \ominus & \text{ if } P(y \mid \mathbf{x}) < th \end{aligned} \quad (28)$$

- High threshold: high specificity, low sensitivity
- Low threshold: low specificity, high sensitivity
- We should select the best threshold for the trade-off between the cost of *false positives* vs *false negatives*



ROC Curve

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

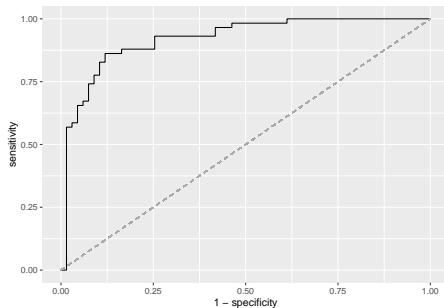
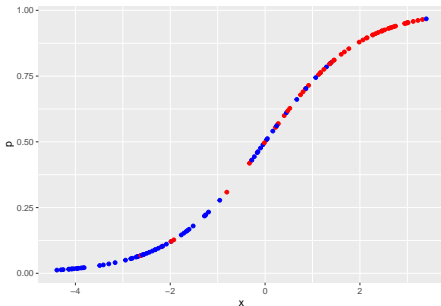
Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- The receiver operating characteristic (*ROC*) curve is plot which shows the performance of a binary classifier as function of its cut-off threshold.
- It essentially shows the *true positive rate (sensitivity)* against the *false positive rate (1-specificity)* for various threshold values.
- The area under the curve (*AUC*) is an aggregated measure of performance.





Multi-class classification problems

- **Email foldering/tagging:** Work (1), Friends (2), Family (3), Hobby (4)
- **Medical diagrams:** Not ill, Cold, Flu
- **Weather:** Sunny, Cloudy, Rain, Snow

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization



Visual of Binary vs Multi-class classification

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

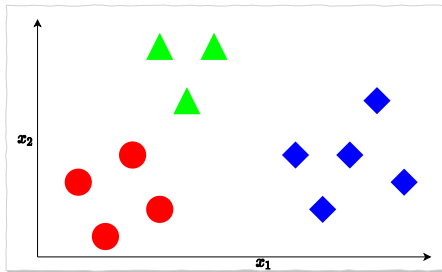
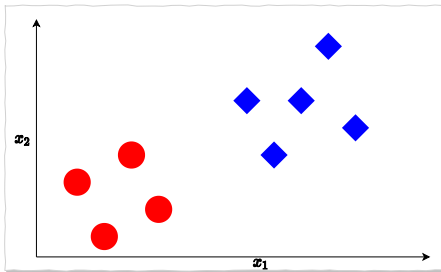
Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization





Approaches

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- One-vs-one
- Hierarchical
- **One-vs-all**



One-vs-all

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

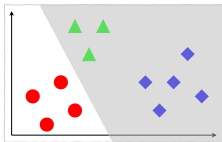
Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

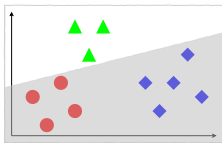
Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization



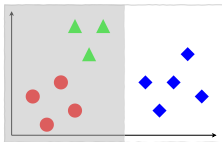
- Class \circ (1):

$$h^{(1)}(\mathbf{x}) = P(y = 1 \mid \mathbf{x}, \mathbf{w}_1)$$



- Class \triangle (2):

$$h^{(2)}(\mathbf{x}) = P(y = 2 \mid \mathbf{x}, \mathbf{w}_2)$$



- Class \diamond (3):

$$h^{(3)}(\mathbf{x}) = P(y = 3 \mid \mathbf{x}, \mathbf{w}_3)$$



One-vs-all (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

Learning

- Train a logistic regression classifier $h^{(i)}(\mathbf{x})$ for each class i to predict the probability that $y = i$.

Prediction

- On a new input \mathbf{x} , to make a prediction, pick the class i that maximizes

$$\arg \max_i (h^{(i)}(\mathbf{x})) \quad (29)$$



Decision boundaries and decision regions

Linear Regression

Simple Linear Model

Weighted Linear Model

Linear Basis Function Model

Classification

Logistic Regression

Binary Classification

Evaluation

Multi-class Classification

Softmax Regression

Softmax Regression

Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

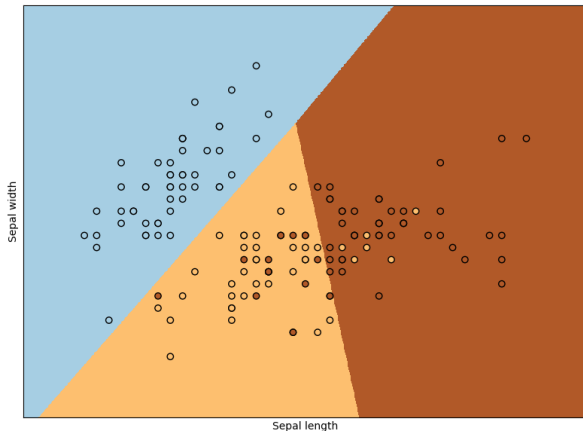
Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- Logistic regression for **Iris** dataset





Softmax Regression

- Softmax Regression
- Cross Entropy vs. MSE



Softmax Regression

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

Concept 4

Softmax regression is a generalization of logistic regression that we can use for multi-class classification

- In Softmax regression, we replace the sigmoid function by the so-called softmax function $\phi(\cdot) = \{\phi_1, \dots, \phi_C\}$.



Score function

Concept 5

The **score function** f that maps the raw features to class scores.

$$\mathbf{z} = f(\mathbf{x}; \mathbf{W}, \mathbf{b}) = \mathbf{W}\mathbf{x} + \mathbf{b} \quad (30)$$

input image



\mathbf{W}

0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

\times

\mathbf{x}

56
231
24
2

$+$

\mathbf{b}

1.1
3.2
-1.2

$=$

\mathbf{z}

-96.8
437.9
61.95

cat score

dog score

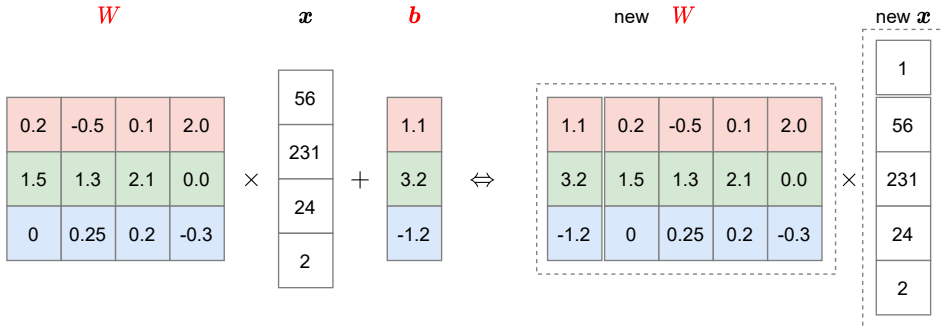
ship score



Score function (cont.)

- Use bias trick $W_0 = \mathbf{b}$ to represent the two parameters W, \mathbf{b} as one

$$\mathbf{z} = f(\mathbf{x}; W) = W\mathbf{x} \quad (31)$$





Softmax function

Concept 6

The **softmax function** converts a score vector $\mathbf{z} = (z_1, \dots, z_C)$ to a discrete distribution vector $\mathbf{p} = (p_1, \dots, p_C)$

$$p_i = P(y = i \mid \mathbf{z}) = \phi_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}, \quad i \in [1, \dots, C] \quad (32)$$

where

$$z_i = w_{i0} + w_{i1}x_1 + \dots + w_{iD}x_D = \mathbf{w}_i^T \mathbf{x} \quad (33)$$



Softmax function (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

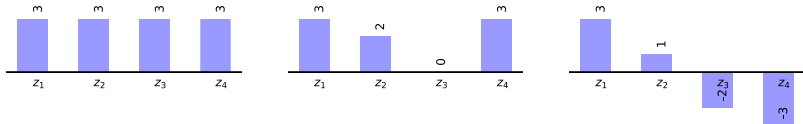
Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

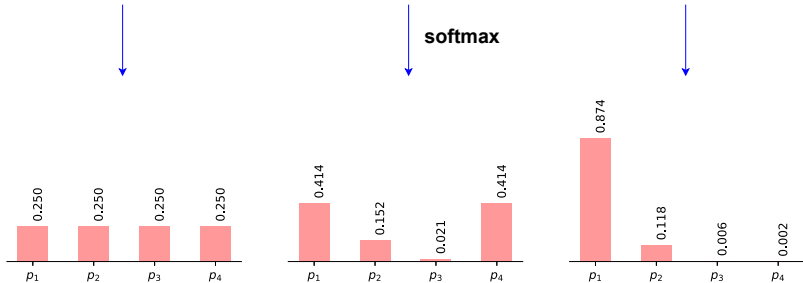
Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

score vectors



softmax



distribution vectors



Softmax function (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

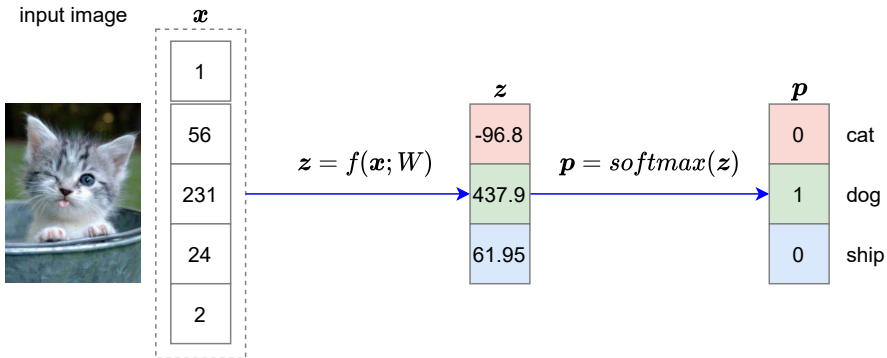
Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization





Problem Statement

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Given $\mathcal{D} = \{(\mathbf{x}_1, y_1) \dots (\mathbf{x}_N, y_N)\}$ where $y_n \in \{1, \dots, C\}$.
- Denote \mathbf{t}_n is a **one-hot encoding** of y_n (or target discrete distribution)
- **Learning goal:** Find a softmax function $\phi_{\mathbf{w}}(\cdot) = \{\phi_1, \dots, \phi_C\}$ that minimize

$$\arg \min_{\mathbf{w}} E(\phi_{\mathbf{w}}) = \arg \min_{\mathbf{w}} \sum_{n=1}^N CE(\mathbf{p}_n, \mathbf{t}_n) \quad (34)$$



Cross Entropy

Concept 7

Cross-entropy (CE) is a measure of the difference between two probability distributions. The cross-entropy between a “true” distribution $\mathbf{t} = (t_1, \dots, t_C)$ and an estimated distribution $\mathbf{p} = (p_1, \dots, p_C)$ is defined as

$$CE(\mathbf{p}, \mathbf{t}) = - \sum_{i=1}^C t_i \log p_i \quad (35)$$



Cross Entropy (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

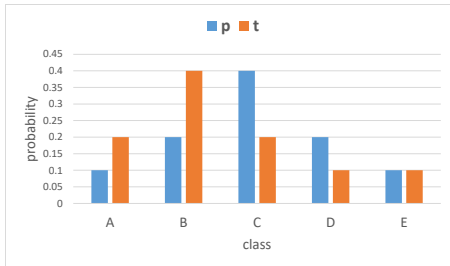
Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization



$$\left. \begin{array}{l} \mathbf{p} = (0.1, 0.2, 0.4, 0.2, 0.1) \\ \mathbf{t} = (0.2, 0.4, 0.2, 0.1, 0.1) \end{array} \right\} \rightarrow CE(\mathbf{p}, \mathbf{t}) = 1.678$$

- $CE > 0$
- $CE(\mathbf{p}, \mathbf{t}) \neq CE(\mathbf{t}, \mathbf{p})$
- CE minimize if $p_i = t_i, \forall i$



MSE

Concept 8

Mean squared error (MSE) is a measure of the average of the squares of the errors.

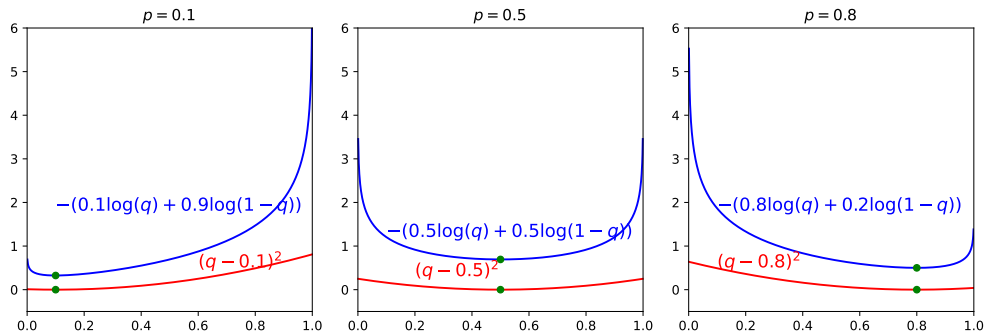
$$MSE(\mathbf{p}, \mathbf{t}) = \frac{1}{C} \sum_{i=1}^C (p_i - t_i)^2 \quad (36)$$

- $MSE \geq 0$
- $MSE(\mathbf{p}, \mathbf{t}) = MSE(\mathbf{t}, \mathbf{p})$
- $MSE = 0$ if $p_i = t_i, \forall i$



Cross Entropy vs. MSE

- Consider three “true” binary distributions $\mathbf{p} = (0.1, 0.9)$, $(0.5, 0.5)$ and $(0.8, 0.2)$





Learning Algorithm

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

1. Initialize the weights W_0 (parameters) at $t = 0$
2. For $t = 1, 2, 3, \dots$ do
 - 2.1 Compute the output distribution p_n for each x_n ($n = 1 \dots N$)

$$p_n = \text{softmax}(W_t x_n) \quad (37)$$

- 2.2 Compute the gradient

$$\nabla_{W_t} E = \frac{1}{N} \sum_{n=1}^N (p_n - t_n) x_n^T \quad (38)$$

- 2.3 Update the weights

$$W_{t+1} = W_t - \eta \nabla_{W_t} E \quad (39)$$

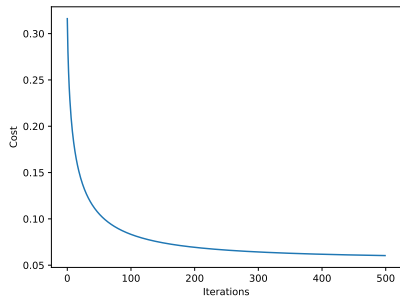
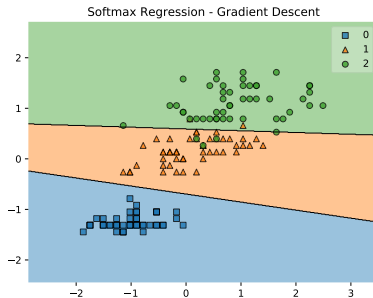
Iterate the next step until W is not change

3. Return the final weights W

Example



- Softmax regression for **Iris** dataset





Capacity, Overfitting and Underfitting

- Model Capacity
- Model vs. Data
- Bias-Variance
- Tradeoff of Capacity
- Regularization
- Tradeoff of Regularization



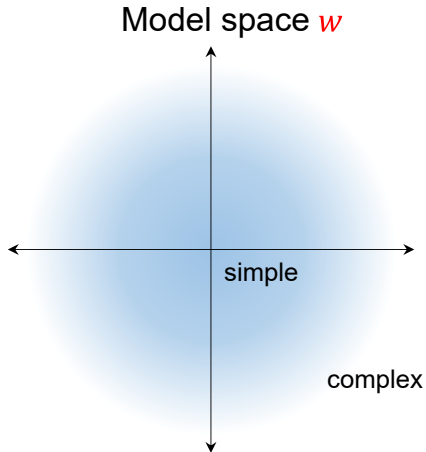
Model Capacity

Concept 9

Capacity is model complexity.

The most common ways to estimate the capacity of a model:

- VC dimension
- The number of parameters
- The norm of parameters





Model vs. Data

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

Concept 10

Data is divided into three sets: **training set**, **validation set** and **test set**.

Concept 11

Models can be too limited. We can't find a function that fits the data well. This is called **underfitting**.

Concept 12

Models can also be too rich. We don't just model the data, but also the underlying noise. This is called **overfitting**.



Model vs. Data (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

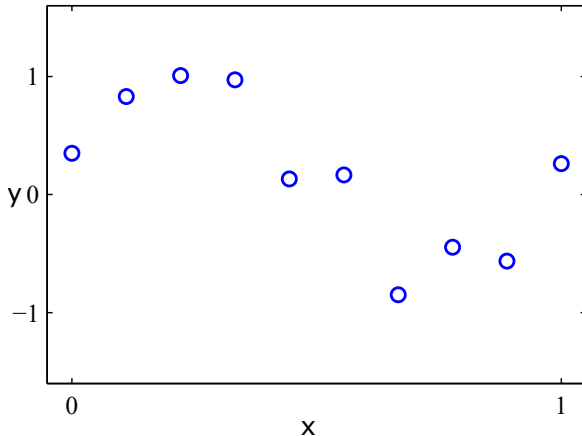
Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Given the data set $\mathcal{D} = \{(x_1, y_1), \dots, (x_{10}, y_{10})\}$ shown in the following figure, find the best regression function to the data





Model vs. Data (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

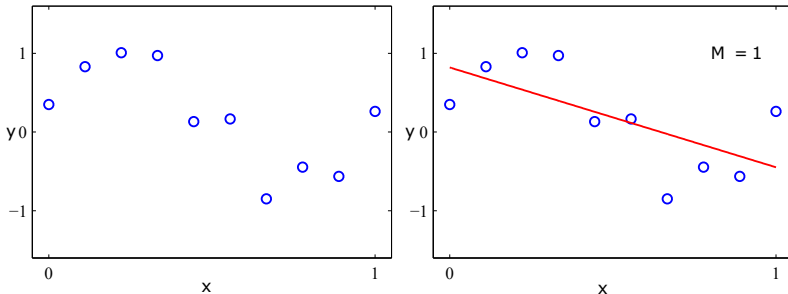
Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Consider the simple hypothesis set

$$\mathcal{H}_1 = \{h \mid y = h(x) = w_0 + w_1 x\} \quad (40)$$



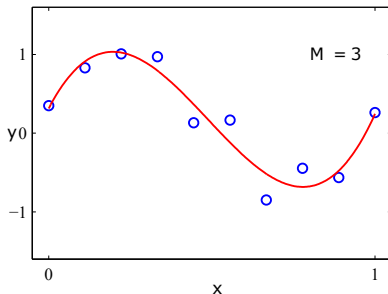
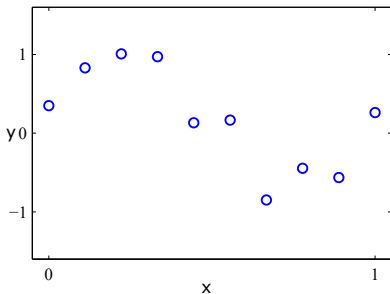
Model vs. Data (cont.)



- Consider the hypothesis set

$$\mathcal{H}_3 = \{h \mid y = h(x) = w_0 + w_1x + w_2x^2 + w_3x^3\} \quad (41)$$

(note that $\mathcal{H}_1 \subset \mathcal{H}_3$)





Model vs. Data (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

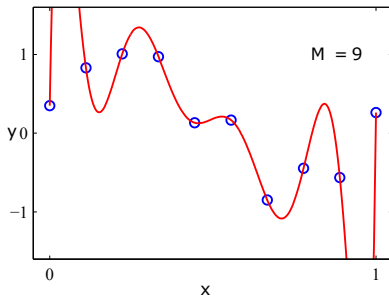
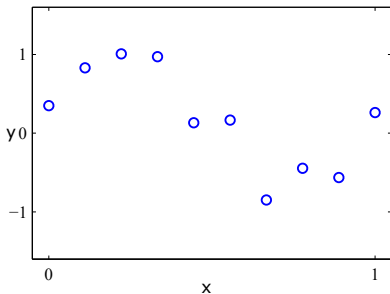
Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- Consider the hypothesis set

$$\mathcal{H}_9 = \{h \mid y = h(x) = w_0 + w_1x + w_2x^2 + \dots + w_9x^9\} \quad (42)$$

(note that $\mathcal{H}_1 \subset \mathcal{H}_3 \subset \mathcal{H}_9$)





Model vs. Data (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

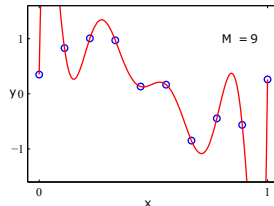
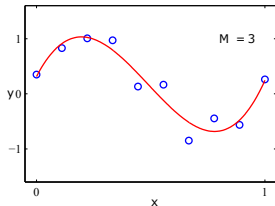
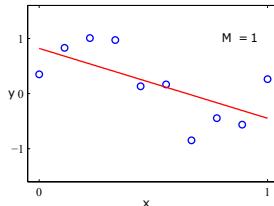
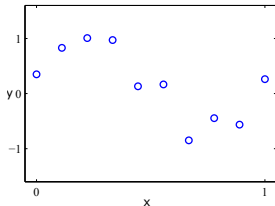
Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

Which one

- Under-fitting
- Over-fitting
- **Appropriate fitting**



Model vs. Data (cont.)



	$M = 1$	$M = 3$	$M = 9$
w_0	0.82	0.31	0.35
w_1	-1.27	7.99	232.37
w_2		-25.43	-5321.83
w_3		17.37	48568.31
w_4			-231639.30
w_5			640042.26
w_6			-1061800.52
w_7			1042400.18
w_8			-557682.99
w_9			125201.43



Model Performance

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity

Model vs. Data

Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

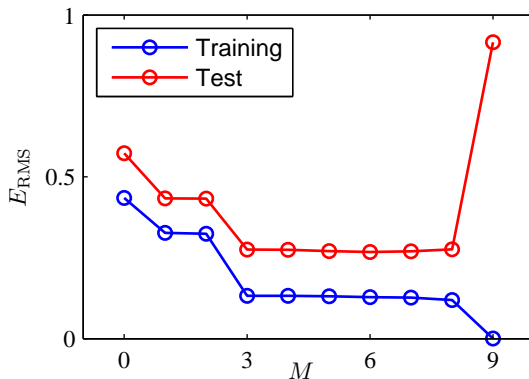


Figure 1: Graphs of the root-mean-square error evaluated on the **training set** and on an **independent test set** for various values of M



What happen if increasing N

Linear Regression

- Simple Linear Model
- Weighted Linear Model
- Linear Basis Function Model

Classification

Logistic Regression

- Binary Classification
- Evaluation
- Multi-class Classification

Softmax Regression

- Softmax Regression
- Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

- Model Capacity
- Model vs. Data
- Bias-Variance
- Tradeoff of Capacity
- Regularization
- Tradeoff of Regularization

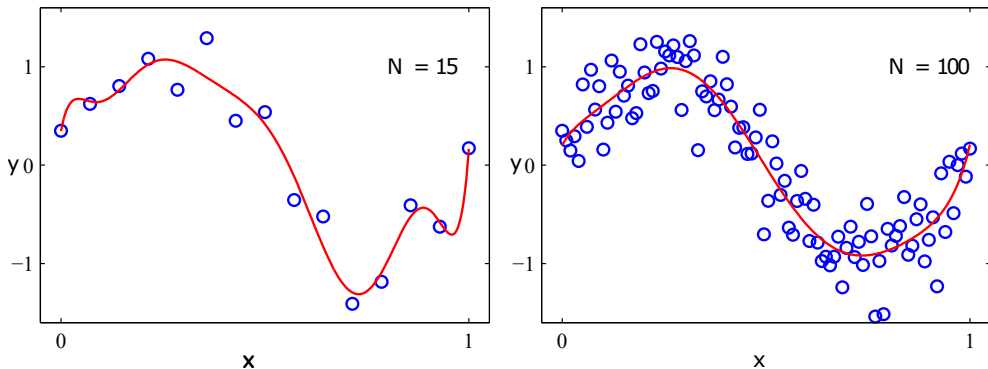


Figure 2: Using the $M = 9$ polynomial for $N = 15$ data points (left plot) and $N = 100$ data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem



Errors in Learning Model

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data

Bias-Variance

Tradeoff of Capacity
Regularization
Tradeoff of Regularization

- **Bias errors:** error due to assumption in the model
 - High bias to signify underfitting
- **Variance errors:** It measures the variability in the results given by model when the dataset is changed
 - High variance to signify overfitting

$$\text{Expected error} = \text{Bias} + \text{Variance} + \text{Irreducible Error} \quad (43)$$



Errors in Learning Model (cont.)

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data

Bias-Variance

Tradeoff of Capacity
Regularization
Tradeoff of Regularization

Concept 13

Given a learning model $\langle \mathcal{H}, \mathcal{A} \rangle$, we define the “average” hypothesis $\bar{g}(x)$

$$\bar{g}(x) = \mathbb{E}_{\mathcal{D}} [g_{\mathcal{D}}(x)] \quad (44)$$

where $g_{\mathcal{D}}(x)$ is the “**best**” hypothesis given the data set \mathcal{D}

- Bias of learning model

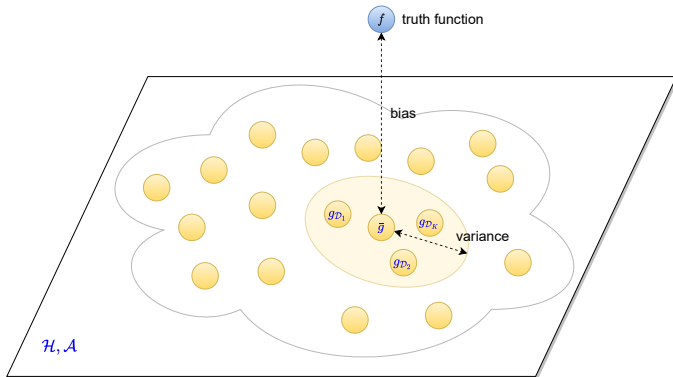
$$\text{Bias} = \mathbb{E}_x \left[(\bar{g}(x) - f(x))^2 \right] \quad (45)$$

where $f(x)$ is the “**truth**” function

- Variance of learning model

$$\text{Variance} = \mathbb{E}_x \left[\mathbb{E}_{\mathcal{D}} \left[(g_{\mathcal{D}}(x) - \bar{g}(x))^2 \right] \right] \quad (46)$$

Errors in Learning Model (cont.)



- Given many independent data sets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$, we can estimate $\bar{g}(x)$ by

$$\bar{g}(x) \approx \frac{1}{K} \sum_{i=1}^K g_{\mathcal{D}_i}(x) \quad (47)$$



Example: two learning models

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance

Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- Consider a target function sine

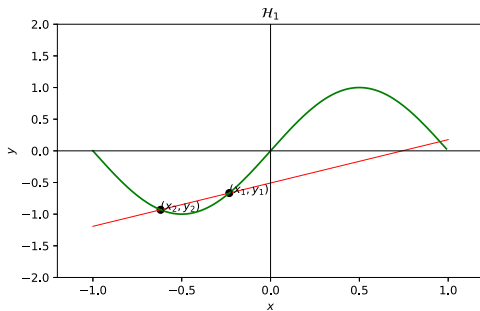
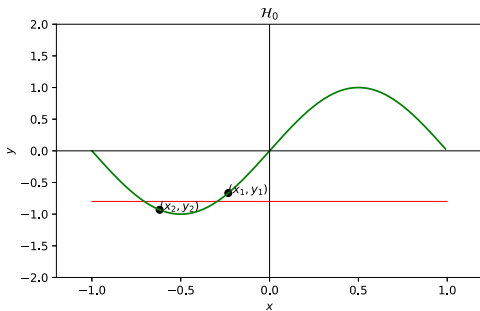
$$\begin{aligned} f &: [-1, 1] \rightarrow \mathbb{R} \\ x &\rightarrow \sin(\pi x) \end{aligned} \quad (48)$$

- We generate 100 data sets $\{\mathcal{D}_i\}$, $i = 1, \dots, 100$, each containing $N = 2$ data points, independently from the sinusoidal curve $f(x) = \sin(\pi x)$. For each data set \mathcal{D}_i , we fit the data using one of two models
 - \mathcal{H}_0 : set of all lines of the form $h(x) = b$
 - \mathcal{H}_1 : set of all lines of the form $h(x) = ax + b$
 - Note that $\mathcal{H}_0 \subset \mathcal{H}_1$

Example: two learning models (cont.)



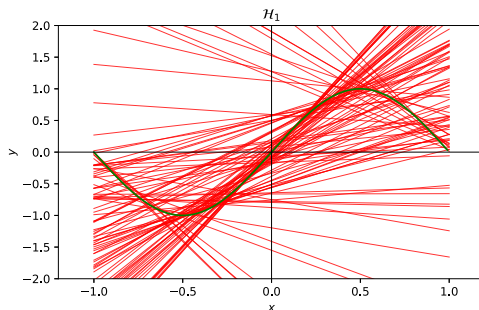
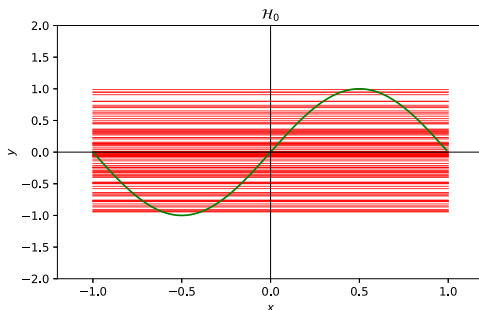
- Given a data set $\mathcal{D}_i = \{(x_1, y_1), (x_2, y_2)\}$.
 - For \mathcal{H}_0 , we choose the constant hypothesis that best fits the data (the horizontal line at the midpoint, $b = (y_1 + y_2)/2$).
 - For \mathcal{H}_1 , we choose the line that passes through the two data points (x_1, y_1) and (x_2, y_2) .



Example: two learning models (cont.)



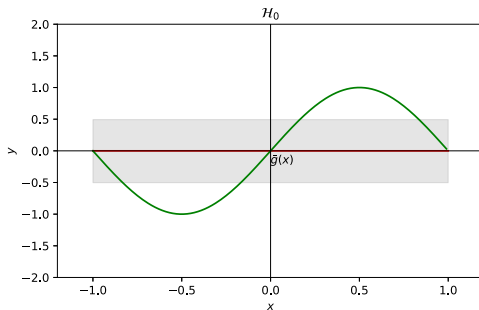
- Repeating this process with 100 data sets $\{\mathcal{D}_i\}, i = 1, \dots, 100$,



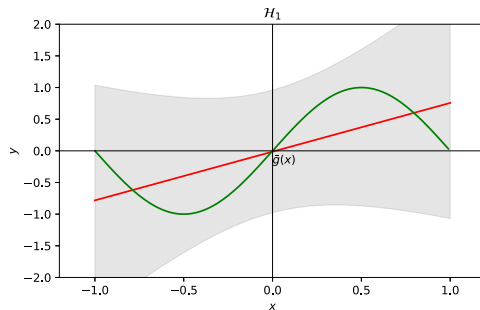
Example: two learning models (cont.)



- The bias-variance for each learning model



- bias=0.50
- var=0.25



- bias=0.21
- var=1.69



Generalization and Capacity

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

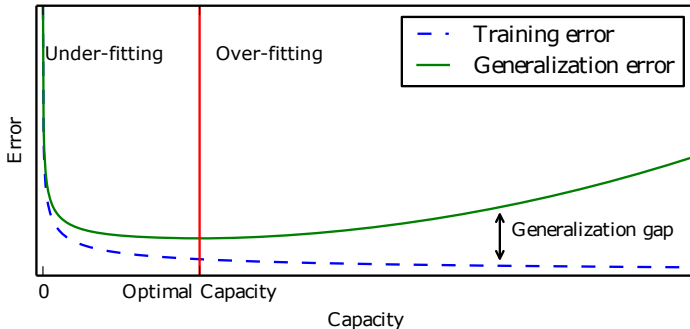
Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization
Tradeoff of Regularization

The criteria determining how well a machine learning model will perform:

1. Make the training error small.
2. Make the gap between training and test (generalization) error small.

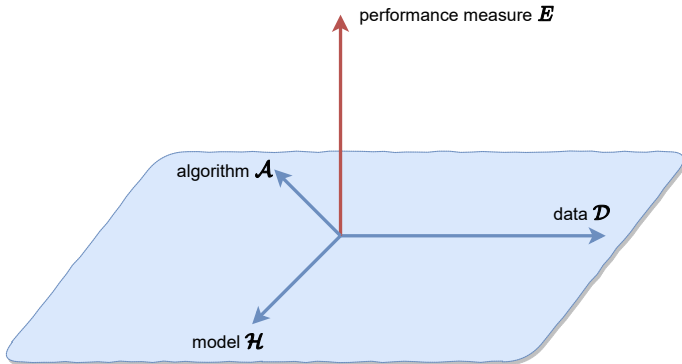




Regularization

Concept 14

Regularization is any modification we make to a learning model that is intended to reduce its generalization error but not its training error.





Addressing

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity

Regularization

Tradeoff of Regularization

High bias	High variance
Obtain more features	Decrease number of features
Decrease regularization λ	Increase regularization λ
Extend model	Obtain more data
Train longer	Stop early
New model architecture	New model architecture



Regularization for Linear Regression

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity

Regularization

Tradeoff of Regularization

- Using regularized MSE_{train} for linear regression

$$\mathbf{w} = \arg \min_{\mathbf{w}} \left[MSE_{train} + \lambda \left(\frac{1}{N} \mathbf{w}^T \mathbf{w} \right) \right] \quad (49)$$

where λ is the regularization coefficient (*hyper-parameter*) that controls the relative importance of the data-dependent error MSE_{train} and the regularization term $\lambda \left(\frac{1}{N} \mathbf{w}^T \mathbf{w} \right)$

- Solving for \mathbf{w} , we obtain

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (50)$$



Example: one learning model

Linear Regression

Simple Linear Model
Weighted Linear Model
Linear Basis Function Model

Classification

Logistic Regression

Binary Classification
Evaluation
Multi-class Classification

Softmax Regression

Softmax Regression
Cross Entropy vs. MSE

Capacity, Overfitting and Underfitting

Model Capacity
Model vs. Data
Bias-Variance
Tradeoff of Capacity
Regularization

Tradeoff of Regularization

- Consider a target function sine

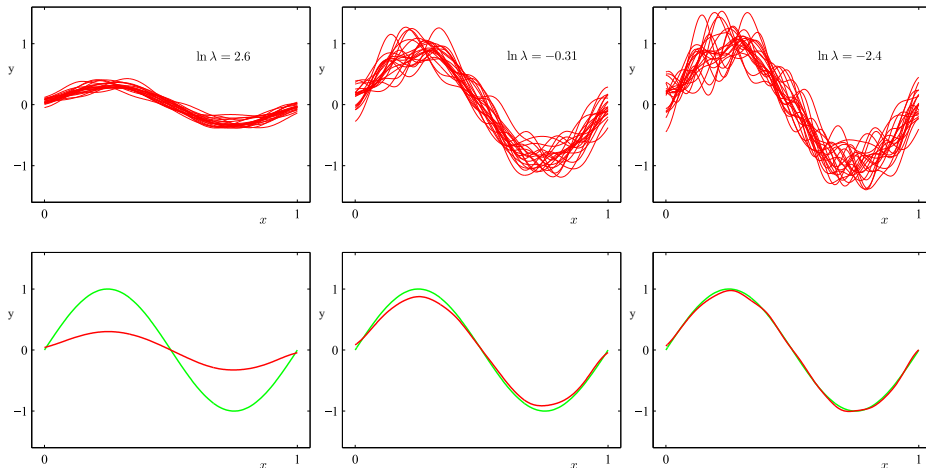
$$\begin{aligned} f &: [-1, 1] \rightarrow \mathbb{R} \\ x &\rightarrow \sin(2\pi x) \end{aligned} \quad (51)$$

- We generate 100 data sets $\{\mathcal{D}_i\}, i = 1, \dots, 100$, each containing $N = 25$ data points, independently from the sinusoidal curve $f(x) = \sin(2\pi x)$. For each data set \mathcal{D}_i , we fit a model with 24 Gaussian basis functions by minimizing the regularized error function

Example: one learning model (cont.)



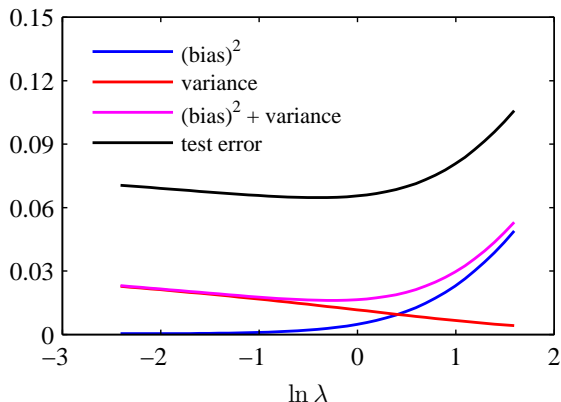
- Illustration of the dependence of bias and variance on model regularization coefficient



Example: one learning model (cont.)



- Summary of the dependence of bias and variance on model regularization coefficient



References



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

Deep learning.

MIT press.



Lê, B. and Tô, V. (2014).

Cở sở trí tuệ nhân tạo.

Nhà xuất bản Khoa học và Kỹ thuật.



Russell, S. and Norvig, P. (2021).

Artificial intelligence: a modern approach.

Pearson Education Limited.