# UVA CS 6316:
# Machine Learning

# Lecture 14: Logistic Regression

Dr. Yanjun Qi

University of Virginia

Department of Computer Science

# Course Content Plan ➜
## Six major sections of this course

❑ ~~Regression (supervised)~~ ← Y is a continuous

❑ Classification (supervised) ← Y is a discrete
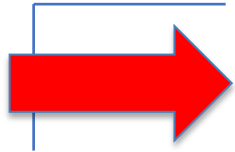
❑ Unsupervised models ← NO Y

❑ Learning theory ← About f()

❑ Graphical models ← About interactions among X1,... Xp
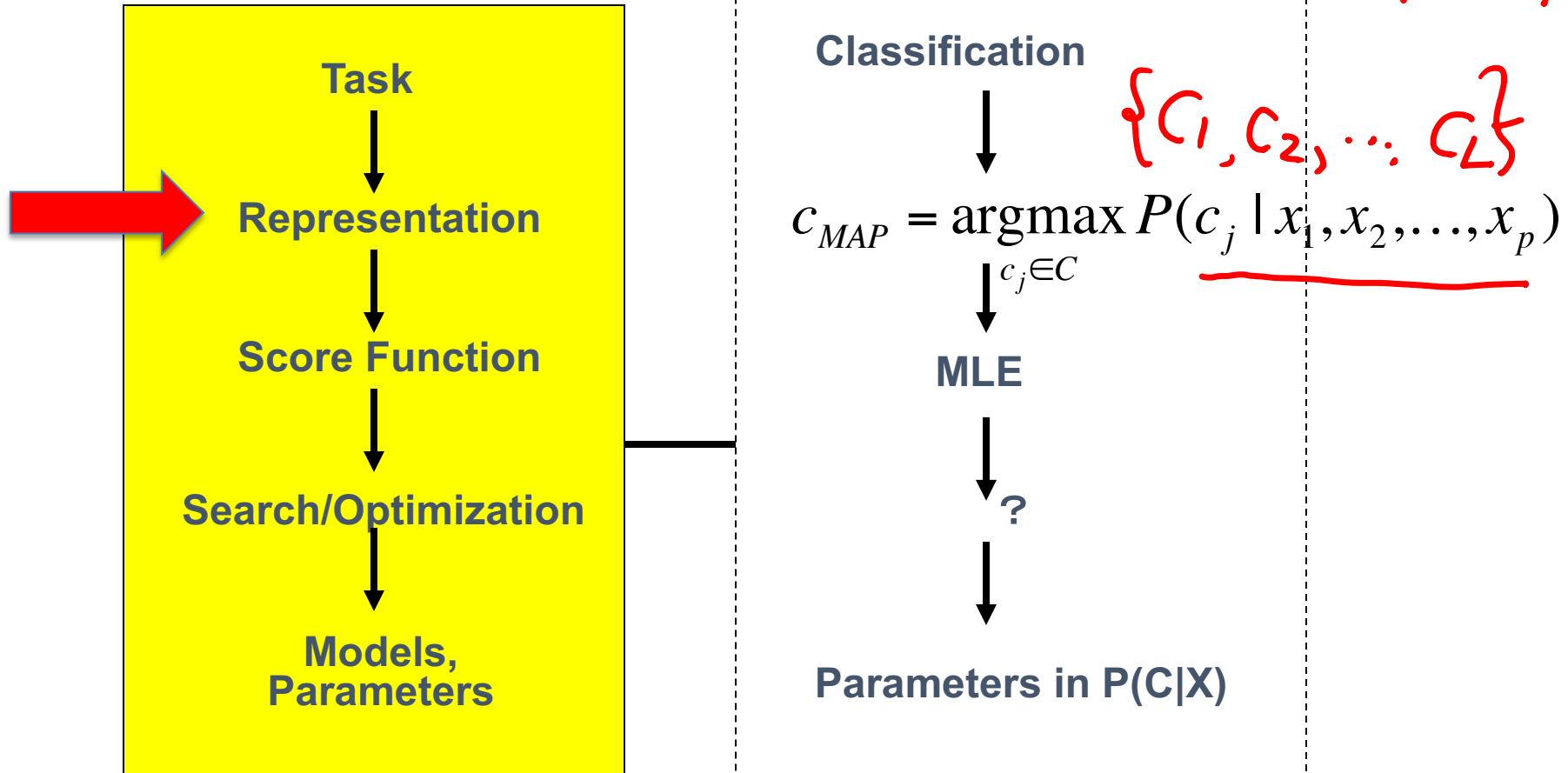
❑ Reinforcement Learning ← Learn program to Interact with its environment

# **Today**

- ❑ Bayes Classifier
- ❑ Logistic Regression
- ❑ Training LG by MLE

# Bayes Classifier

$c^* = \text{argmax } P(c_j \mid x_1, \ldots x_p)$

$\{c_1, c_2, \ldots, c_L\}$

**Classification**

$$c_{MAP} = \underset{c_j \in C}{\text{argmax}} \, P(c_j \mid x_1, x_2, \ldots, x_p)$$

**MLE**

**?**

**Parameters in P(C|X)**

**Task**

**Representation**

**Score Function**

**Search/Optimization**

**Models, Parameters**

# Bayes classifiers

- Treat each feature attribute and the class label as random variables.

$$\{c_1, \cdots c_L\}$$

# Bayes classifiers

- Treat each feature attribute and the class label as random variables.

- Testing: Given a sample **x** with attributes ( $x_1, x_2, \ldots, x_p$ ):
  - Goal is to predict its class $c$.
  - Specifically, we want to find the class that maximizes $p( c \mid x_1, x_2, \ldots, x_p )$.

- Training: can we estimate $p(C_i \mid \mathbf{x}) = p( C_i \mid x_1, x_2, \ldots, x_p )$ directly from data?

# Bayes Classifiers – MAP Rule

*Task*: Classify a new instance *X* based on a tuple of attribute values $X = \langle X_1, X_2, \ldots, X_p \rangle$ into one of the classes

$$c_{MAP} = \underset{c_j \in C}{\operatorname{argmax}} \, P(c_j \mid x_1, x_2, \ldots, x_p)$$

MAP Rule

MAP = Maximum Aposteriori Probability

Adapt From Carols' prob tutorial

# Bayes Classifiers – MAP Classification Rule

- Establishing a probabilistic model for classification
  - ➔ **MAP** classification rule
    - **MAP**: **M**aximum **A P**osterior
    - Assign $x$ to $c^*$ if

$$\sum_{j=1}^{L} P(C = c_j | x) = 1$$

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x})$$

$$\text{for } c \neq c^*, \quad c = c_1, \cdots, c_L$$

Adapt from Prof. Ke Chen NB slides

$$f : \boxed{X} \longrightarrow \boxed{C}$$

Output as Discrete Class Label $C_1, C_2, \ldots, C_L$

Establishing a probabilistic model for classification

$$c_{MAP} = \underset{c_j \in C}{\mathrm{argmax}}\, P(c_j \mid x_1, x_2, \ldots, x_p)$$

$$\frac{P(X,c)}{P(X)}$$

**Generative**

$$\underset{c \in C}{\mathrm{argmax}}\, P(c \mid X) = \underset{c \in C}{\mathrm{argmax}}\, P(X,c) = \underset{c \in C}{\mathrm{argmax}}\, P(X \mid c) P(c)$$

Later!

**Discriminative**

$$\underset{c \in C}{\mathrm{arg\,max}}\, P(c / \mathbf{X}) \quad C = \{c_1, \cdots, c_L\}$$

# Recap: Statistical Decision Theory (Extra)

- Random input vector: X
- Random output variable: Y
- Joint distribution: Pr(X,Y )  $\Rightarrow D = \begin{pmatrix} (\vec{x}_1, y_1) \\ \vdots \\ (\vec{x}_n, y_n) \end{pmatrix}$
- Loss function L(Y, f(X))
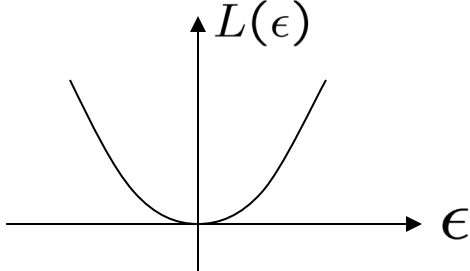
- Expected prediction error (EPE):

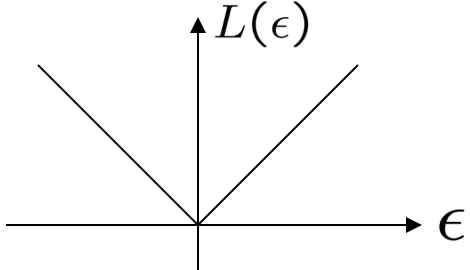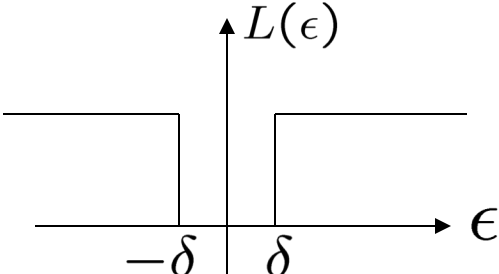$$\mathrm{EPE}(f) = \mathrm{E}(L(Y, f(X))) = \int L(y, f(x)) \Pr(dx, dy)$$

$$\mathrm{e.g.} = \int (y - f(x))^2 \Pr(dx, dy)$$

e.g. Squared error loss (also called L2 loss )

Consider population distribution

# SUMMARY: WHEN Expected prediction error (EPE) USES DIFFERENT LOSS

| Loss Function | Estimator $\hat{f}(x)$ |
|---|---|
| $L_2$ <br>  | $EPE = E_{X,Y}\left(Y - f(x)\right)^2$ <br> $\hat{f}(x) = E[Y\|X = x]$ |
| $L_1$ <br>  | $\hat{f}(x) = \text{median}(Y\|X = x)$ |
| 0-1 <br>  | $\hat{f}(x) = \arg\max_{Y} P(Y\|X = x)$ <br><br> (Bayes classifier / MAP) |

Please read extra slides for WHY MAP-rule makes sense

$$EPE(f) = E_{X,C}\left(L(C, f(X))\right)$$

$$= E_X E_{C|X}\left[L(C, f(X))\,\middle|\, X\right]$$

Discrete RV's Expectation

$$E_C(C) = \sum_{i=1}^{L} c_i p(c_i)$$

$$= E_X \sum_{k=1}^{L} L\left[C_k, f(X)\right] Pr(C_k | X)$$

$$\underset{f}{argmin}\ EPE(f(X))$$

$\Rightarrow$ Pointwise minization     When $X = x$

$$\Rightarrow \hat{f}(X = x) = \underset{f(x) \in C}{argmin} \sum_{k=1}^{L} L(C_k, f(x))\, Pr(C_k | X = x)$$

$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$

$\Downarrow$

$$\Rightarrow \hat{f}(x) = \underset{C_k \in C_1}{argmax}\ Pr(C_k | X = x)$$

$\left\{ \begin{array}{l} C_2 \\ C_3 \\ \vdots \\ C_L \end{array} \right.$

$\left\{ \begin{array}{l} p(C_1 | x) \\ p(C_2 | x) \\ \vdots \\ p(C_L | x) \end{array} \right.$

# Today:

$$X \longrightarrow C \quad : \quad \underbrace{P(C|X)}_{f(X)}$$

 – **Discriminative model**

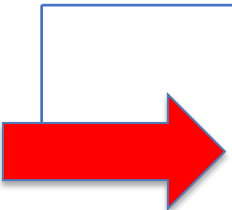$$\arg\max_{c \in C} P(c / \mathbf{X}), \quad C = \{c_1, \cdots, c_L\}$$

$P(c_1 | \mathbf{x}) \quad P(c_2 | \mathbf{x}) \qquad P(c_L | \mathbf{x})$

• • •

## Discriminative Probabilistic Classifier

$x_1 \quad x_2 \qquad \qquad x_p$

• • •

$$\mathbf{x} = (x_1, x_2, \cdots, x_p)$$

# Today

❑ Bayes Classifier

➡ ❑ Logistic Regression

❑ Training LG by MLE

# Logistic Regression



**Task**

**Representation**

**Score Function**

**Search/Optimization**

**Models, Parameters**

**Binary Classification**

**Log-odds = linear function of X's**

**MLE**

**Iterative (Newton) method**

**Logistic weights**

$$P(y=0|x) = 1 - P(y=1|x)$$

$$P(y=1|x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

# Multivariate linear regression to Logistic Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

Logistic regression for binary classification

$$\ln\left[\frac{P(\hat{y}|x)}{1-P(y|x)}\right] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

$$\log\left(\frac{P(y=1|x)}{P(y=0|x)}\right)$$

# Logistic Regression p(y|x)

$$\ln\left[\frac{P(y|x)}{1-P(y|x)}\right] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

$$P(y|x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p}} = \frac{1}{1 + e^{-(\beta_0 + \beta^T X)}}$$

# The logit function View (e.g. when with 1D x)

$$P(y|x) = \frac{e^{\alpha+\beta x}}{1 + e^{\alpha+\beta x}}$$

logisitic $\int$

$$\ln\left(\frac{P = P(y=1|x)}{1-P}\right)$$

$$\ln\left[\frac{P(y|x)}{1 - P(y|x)}\right] = \alpha + \beta x$$

logit / log-odd

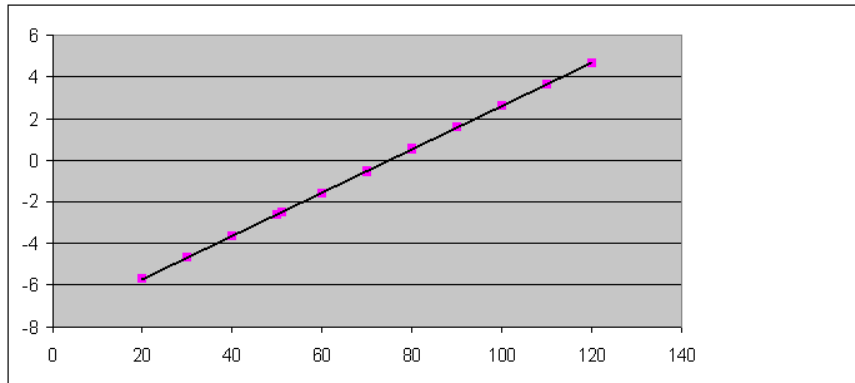Logit function

Logit of $P(y|x)$

# Binary Logistic Regression (Two Views)

ln[p/(1-p)]



x

P (Y=1|x)



$\int$

shape

x

Bernoulli Distribution

$P_{HEAD}$

$\Downarrow$

$Y \in \{ \underset{0,}{T,} \underset{1}{H} \}$

$P_{HEAD} = P(y=1|x)$

$= \dfrac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$

1

0

$P(y=1|x)$   $1-p(y=1x)$

# View I: logit of p(y=1|x) is linear function of x

e.g. Probability of disease

P (Y=1|X)

$$P(y|x) = \frac{e^{\alpha+\beta x}}{1+e^{\alpha+\beta x}}$$



$x$

# View II: "S" shape function compress output to [0,1]

e.g. Probability of disease

P (Y=1|X)

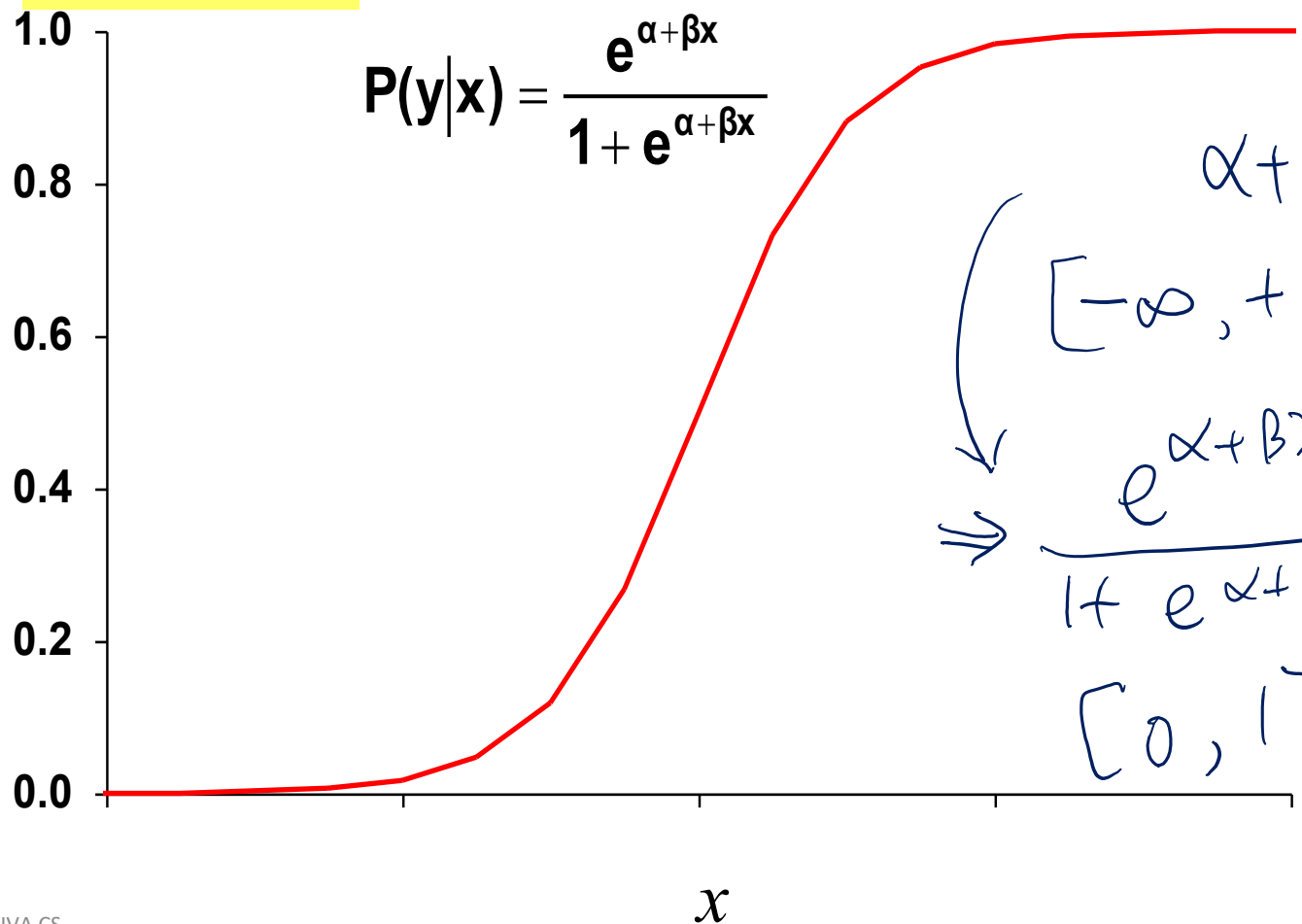$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

1.0

0.8

0.6

0.4

0.2

0.0

$x$

$$\alpha + \beta x$$

$$[-\infty, +\infty]$$

$$\Rightarrow \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

$$[0, 1]$$

# View III: Logistic Regression models a linear classification boundary!



$$y = \{ \underset{1}{H} , \underset{0}{T} \}$$

$$\underset{y \in \{0,1\}}{\text{argmax}} \; P(Y \mid \underline{x})$$

$$\frac{P(y=0 \mid x) = P(y=1 \mid x)}{\log \left( \frac{P(y=1 \mid x)}{P(y=0 \mid x)} \right) = \vec{\beta}^T \vec{x} = \log(1) = 0} \Rightarrow \begin{array}{c} \text{Decision Boundary} \\ \frac{P(y=1 \mid x)}{P(y=0 \mid x)} = 1 \end{array}$$

# Logistic Regression models a linear classification boundary!

$$y \in \{0,1\}$$

$$\ln\left[\frac{P(y|x)}{1 - P(y|x)}\right] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

**Decision Boundary ➜ equals to zero**

$$\ln\left[\frac{P(y=1|x)}{P(y=0|x)}\right] = \ln\left[\frac{P(y=1|x)}{1 - P(y=1|x)}\right] = \alpha + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

# Logistic Regression models a linear classification boundary!

[Seperate two classes]

$$\ln \frac{P(y=1|x)}{1 - P(y=1|x)} = \ln \frac{P(y=1|x)}{P(y=0|x)} = 0$$

0.5

0.5

linear hyperplane

$$\alpha + \beta_1 x_1 + \cdots + \beta_p x_p = 0$$

$$\Uparrow$$

Boundary points

$$p(y=1|x) = p(y=0|x)$$

# Logistic Regression—when?

$\Rightarrow y$ is model with Bernoulli $(p)$

Logistic regression models are appropriate when the target variable is coded as 0/1.

$\Rightarrow p$ is a func of $x$

We only observe "0" and "1" for the target variable— but we think of the target variable conceptually as a probability that "1" will occur.

This means we use Bernoulli distribution to model the target variable with its Bernoulli parameter $p = p(y=1|x)$ predefined.

The main interest ➔ predicting the probability that an event occurs (i.e., the probability that $p(y=1|x)$ ).

# Logistic Regression Assumptions

- Linearity in the logit – the regression equation should have a linear relationship with the logit form of the target variable

- There is no assumption about the feature variables / target predictors being linearly related to each other.
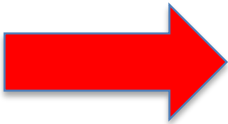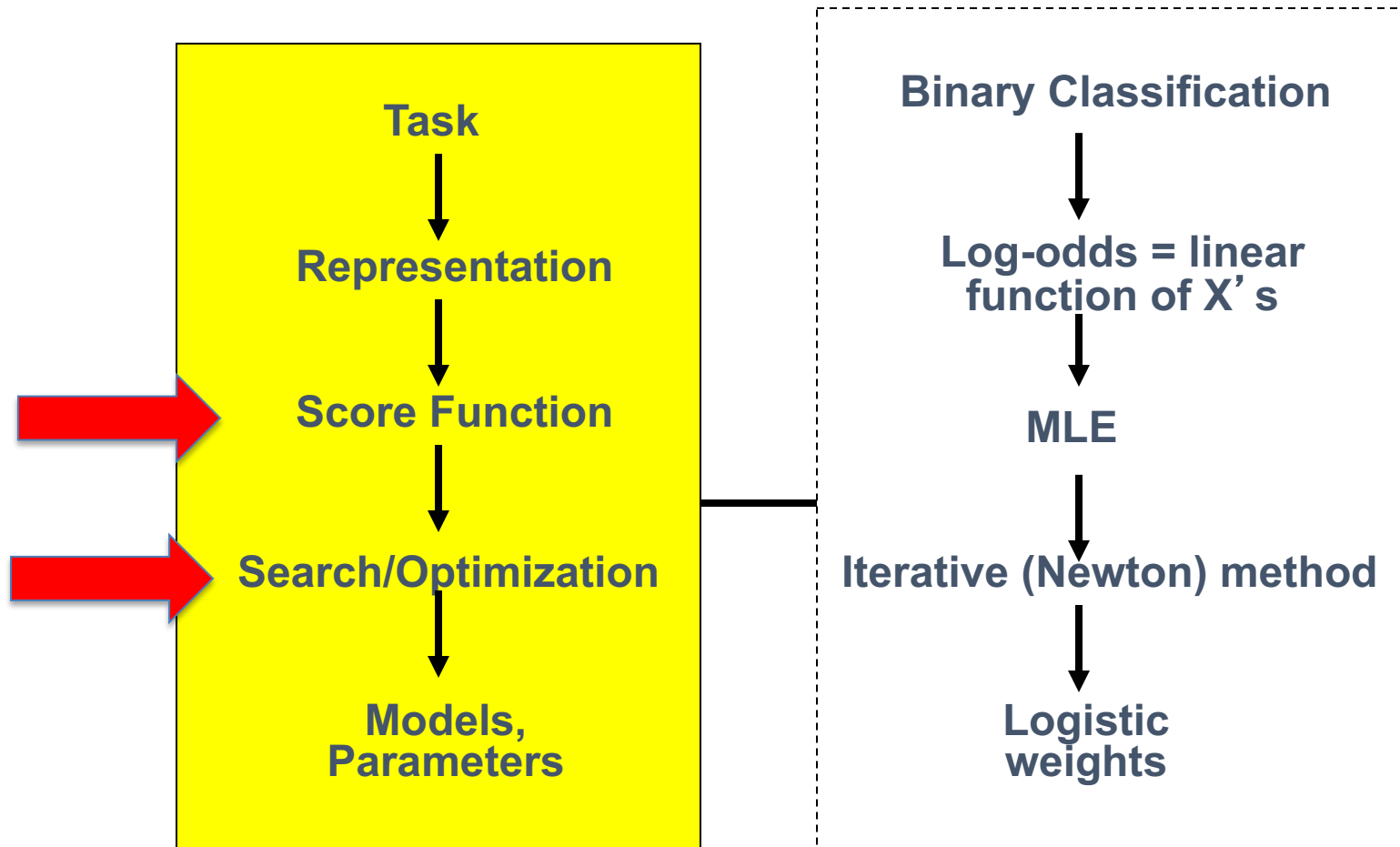


P(y=1|x)     1-p(y=1x)

func of x
with parameter $\vec{\beta}$ to learn from training data

# Today

☐ Bayes Classifier

☐ Logistic Regression

☐ Training LG by MLE

# Logistic Regression



| Task | Binary Classification |
|------|----------------------|
| ↓ | ↓ |
| **Representation** | **Log-odds = linear function of X's** |
| ↓ | ↓ |
| **Score Function** → | **MLE** |
| ↓ | ↓ |
| **Search/Optimization** → | **Iterative (Newton) method** |
| ↓ | ↓ |
| **Models, Parameters** | **Logistic weights** |

$$P(y = 1 | x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

# Review: Maximum Likelihood Estimation

*e.g.*

$$Z = (X_1, \cdots, X_p, Y)$$

*logistic regression*

A general Statement

Consider a sample set T=($Z_1$...$Z_n$) which is drawn from a probability distribution P(Z|\theta) where \theta are parameters.

$$P(Z|\theta)$$

If the Zs are independent with probability density function P($Z_i$|\theta), the joint probability of the whole set is

$$\underset{\theta}{\arg\max} \underbrace{P(Z_1 \ldots Z_n | \theta)}_{\text{data likelihood}} = \prod_{i=1}^{n} P(Z_i | \theta)$$

$$0 < P(z_i | \theta) < 1$$

this may be maximised with respect to \theta to give the maximum likelihood estimates.

# The idea is to

✓ assume a particular model with unknown parameters, $\theta$

The idea is to

✓ assume a particular model with unknown parameters, $\theta$
✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(Z_i|\theta)$

The idea is to

- ✓ assume a particular model with unknown parameters, $\theta$
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(Z_i|\theta)$
- ✓ We have observed a set of outcomes in the real world.

The idea is to

- ✓ assume a particular model with unknown parameters, $\theta$
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters.
  $$P(Z_i|\theta)$$
- ✓ We have observed a set of outcomes in the real world.
  $$Z_1, Z_2, \ldots, Z_n$$
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.
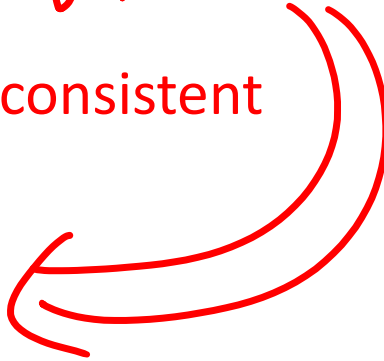
The idea is to

✓ assume a particular model with unknown parameters, $\theta$
✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(Z_i|\theta)$
✓ We have observed a set of outcomes in the real world.
✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \underset{\theta}{argmax}\ P(Z_1 \dots Z_n|\theta) = \prod_{i=1}^{n} P(Z_i|\theta)$$

This is maximum likelihood. In most cases it is both consistent and efficient.

$$log(L(\theta)) = \sum_{i=1}^{n} log(P(Z_i|\theta)$$

It is often convenient to work with the Log of the likelihood function.

The idea is to

- ✓ assume a particular model with unknown parameters, $\theta$
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(Z_i|\theta)$
- ✓ We have observed a set of outcomes in the real world.
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \underset{\theta}{argmax} \ P(Z_1 \ldots Z_n|\theta)$$

Likelihood

This is maximum likelihood. In most cases it is both consistent and efficient.

$$log(L(\theta)) = \sum_{i=1}^{n} log(P(Z_i|\theta))$$

Log-Likelihood

It is often convenient to work with the Log of the likelihood function.

# Review: Defining Likelihood for basic Bernoulli

Given: $\{z_1, z_2, \cdots z_n\}$

$\Downarrow$

$\{H, H, T, \cdots H\}_n$

$\Downarrow$ reformulate

$\{1, 1, 0, \cdots, 1\}_n$

Constant

$\theta = \{p\}$

$= \{p(Head)\}$

$p(z_i | \underline{\theta}) = p^{z_i}(1-p)^{1-z_i}$ (Here $z_i \in \{0,1\}$)

$p(z_i) = \begin{cases} p, & \text{if } z_i = H/1 \\ 1-p, & \text{if } z_i = T/0 \end{cases} \Rightarrow \underset{p}{\text{argmax}} \prod_{i=1}^{n} p^{z_i}(1-p)^{1-z_i}$

# Defining Likelihood

$$\{ H, \ H, \ T, \ \cdots \ H \}_n$$

Observing binary samples $z_i$

Logistic Regression $z_i = y_i | x_i$

$$\{ y_1 | x_1, \ y_2 | x_2, \ \ldots, \ y_n | x_n \}$$

PMF:

$$\Pr(z_i | p) = p^{z_i}(1-p)^{1-z_i}$$

$$P(z_i | \beta)$$

$$= P(y_i | x_i, \beta)$$

LIKELIHOOD:

$$L(p) = \prod_{i=1}^{n} p^{z_i}(1-p)^{1-z_i}$$

↑

function of p=Pr(head)

Now we just rewrite

$$\hat{y}_i = P(y = 1 | x_i)$$

$$P(z_i | \beta) = \hat{y}_i^{y_i}(1-\hat{y}_i)^{1-y_i}$$

LIKELIHOOD:

Basic Bernoulli

$$L(p) = \prod_{i=1}^{n} p^{z_i} (1-p)^{1-z_i}$$

↑

function of p=Pr(head)

Logistic / Bernoulli

$$L(\beta)$$

$$= \prod_{i=1}^{n} P(y_i = 1 | x_i)^{y_i} \left(1 - P(y_i = 1 | x_i)\right)^{1-y_i}$$

$$= \prod_{i=1}^{n} \hat{y}_i^{\,y_i} (1 - \hat{y}_i)^{1-y_i}$$

$$\log(L(p)) = \log\left[\prod_{i=1}^{n} p^{z_i} (1-p)^{1-z_i}\right]$$
$$= \sum_{i=1}^{n} (z_i \log p + (1-z_i)\log(1-p))$$

Loglikehiood

$$\ell\ell(\beta) = \sum_{i=1}^{n} \left( y_i \log \hat{y}_i + (1-y_i)\log(1-\hat{y}_i) \right)$$

$$ll\; l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

When training set includes ($x_i$, $y_i$), i=1,…,$N$

$$ll(\beta) = \sum_{i=1}^{N} \log P(y_i \mid x_i)$$

$$\text{Here } P(y_i \mid x_i) = \begin{cases} P(y=1 \mid x_i), & \text{if } y_i = 1 \\ P(y=0 \mid x_i), & \text{if } y_i = 0 \end{cases}$$

$$= \left(p(y=1 \mid x_i)\right)^{y_i} \left(1 - P(y=1 \mid x_i)\right)^{1-y_i}$$

# MLE for Logistic Regression Training

Training set: $(x_i, y_i)$, i=1,...,$N$

$$l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

$$= \sum_{i=1}^{N} \{y_i \log(\Pr(Y = 1 \mid X = x_i)) + (1 - y_i) \log(\Pr(Y = 0 \mid X = x_i))\}$$

$$= \sum_{i=1}^{N} (y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)}) + (1 - y_i) \log \frac{1}{1 + \exp(\beta^T x_i)})$$

$$= \sum_{i=1}^{N} (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)))$$

Cross entropy loss $\sum_{i=1}^{n} y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)$

# Summary: MLE for Logistic Regression Training

Let's fit the logistic regression model for $K$=2, i.e., number of classes is 2

Training set: $(x_i, y_i)$, i=1,…,$N$

(conditional )
Log-likelihood:

How?

For Bernoulli distribution

$$p(y \mid x)^y (1-p)^{1-y}$$

$$l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

$$= \sum_{i=1}^{N} y_i \log(\Pr(Y = 1 \mid X = x_i)) + (1 - y_i)\log(\Pr(Y = 0 \mid X = x_i))$$

$$= \sum_{i=1}^{N} (y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)}) + (1 - y_i)\log \frac{1}{1 + \exp(\beta^T x_i)})$$

$$= \sum_{i=1}^{N} (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)))$$

$x_i$ are $(p+1)$-dimensional input vector with leading entry 1
\beta is a $(p+1)$-dimensional vector

We want to maximize the log-likelihood in order to estimate \beta

# Logistic Regression

**Task**

↓

**Representation**

↓

**Score Function**

↓

**Search/Optimization**

↓

**Models, Parameters**

**Binary Classification**

↓

**Log-odds = linear function of X's**

↓

**MLE**

↓

**Iterative (Newton) method**

↓

**Logistic weights**

$$P(y = 1 | x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

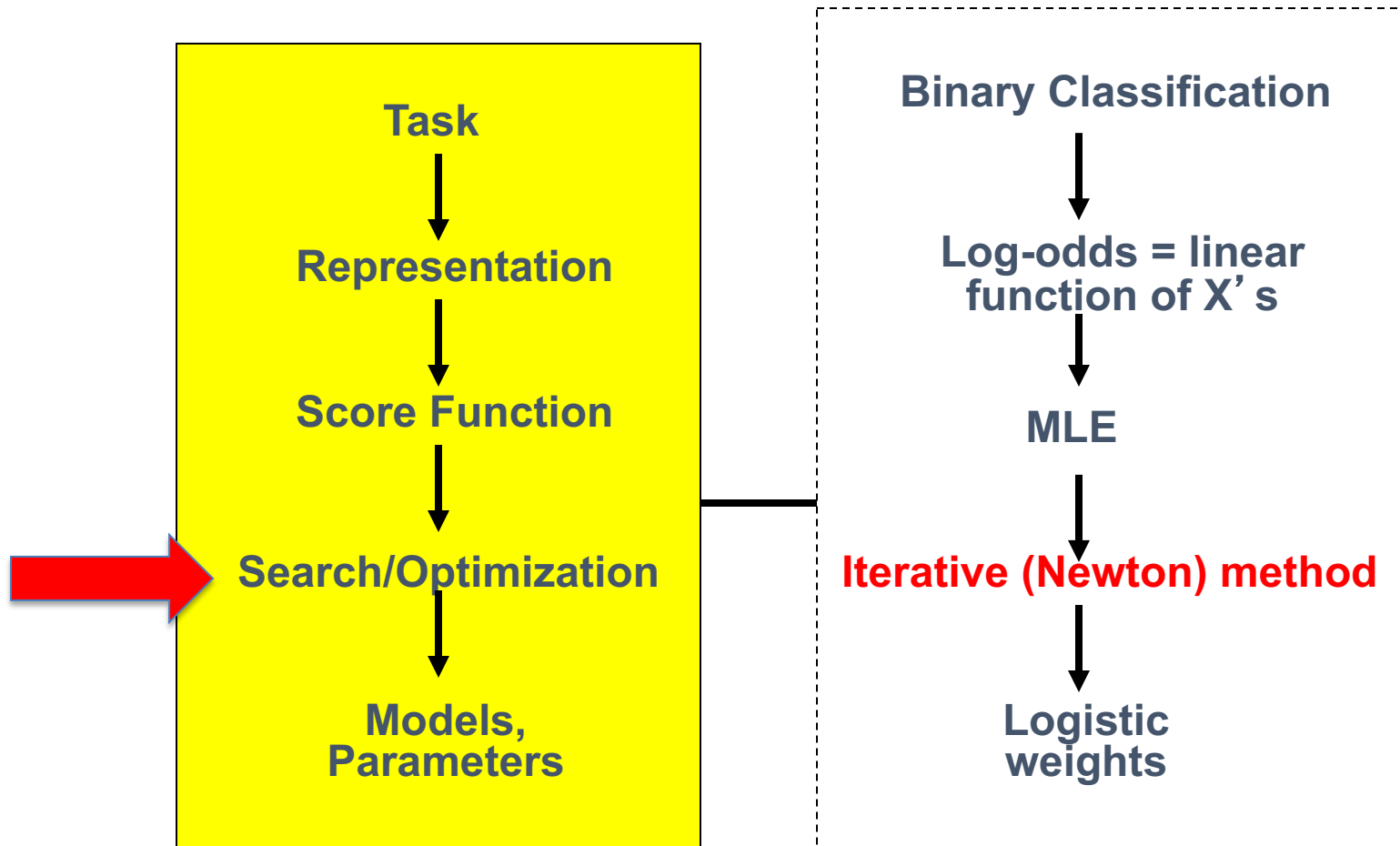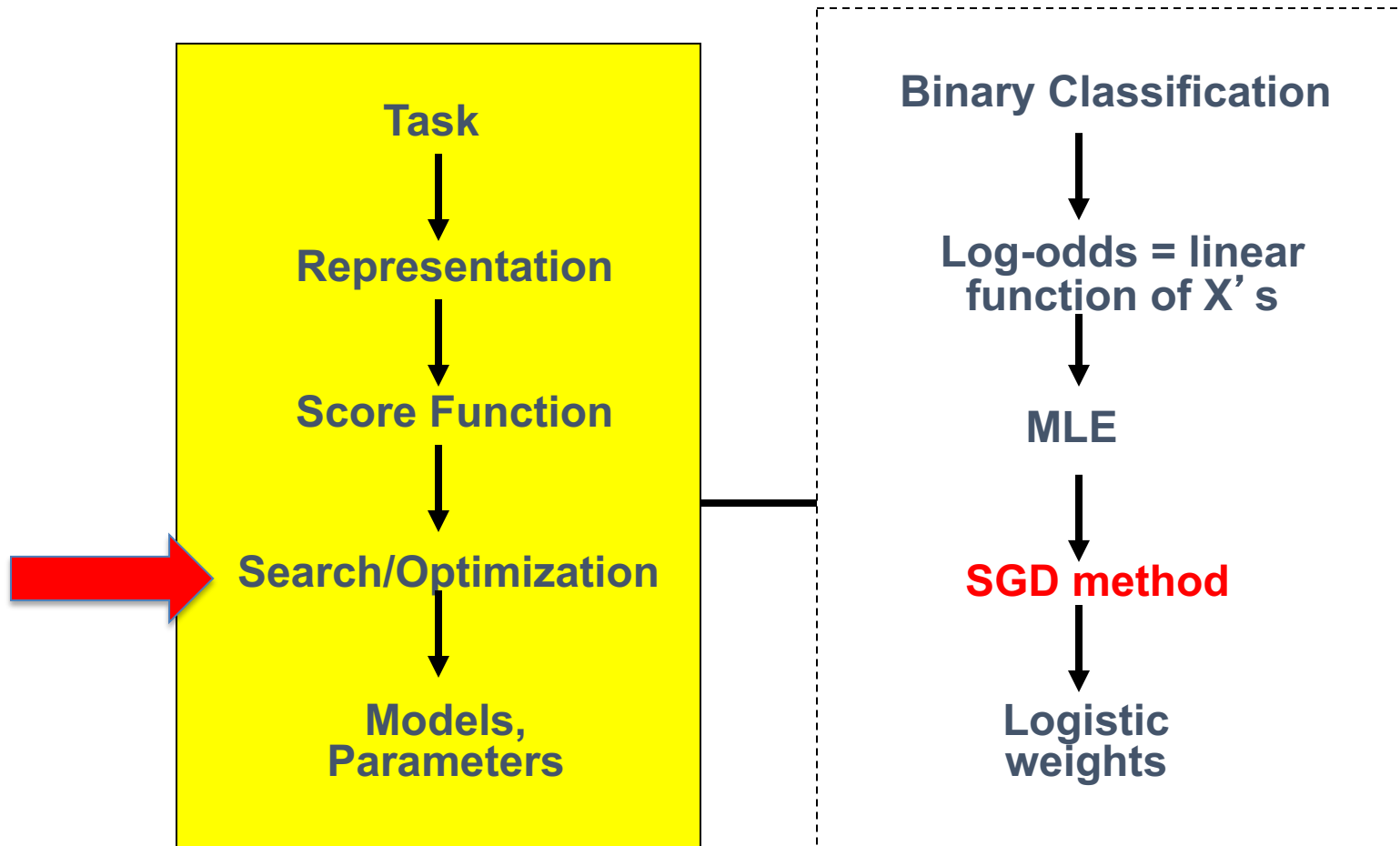# MLE for Logistic Regression Training

Training set: $(x_i, y_i)$, i=1,…,$N$

$$l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

$$= \sum_{i=1}^{N} \{y_i \log(\Pr(Y = 1 \mid X = x_i)) + (1 - y_i)\log(\Pr(Y = 0 \mid X = x_i))\}$$

$$= \sum_{i=1}^{N} \left( y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)} \right) + (1 - y_i)\log \frac{1}{1 + \exp(\beta^T x_i)} )$$

$$= \sum_{i=1}^{N} (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)))$$

See Extra Slides How to used Newton-Raphson optimization

# Logistic Regression



**Task**

↓

**Representation**

↓

**Score Function**

↓

**Search/Optimization**

↓

**Models, Parameters**

**Binary Classification**

↓

**Log-odds = linear function of X's**

↓

**MLE**

↓

**SGD method**

↓

**Logistic weights**

$$P(y = 1 \mid x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

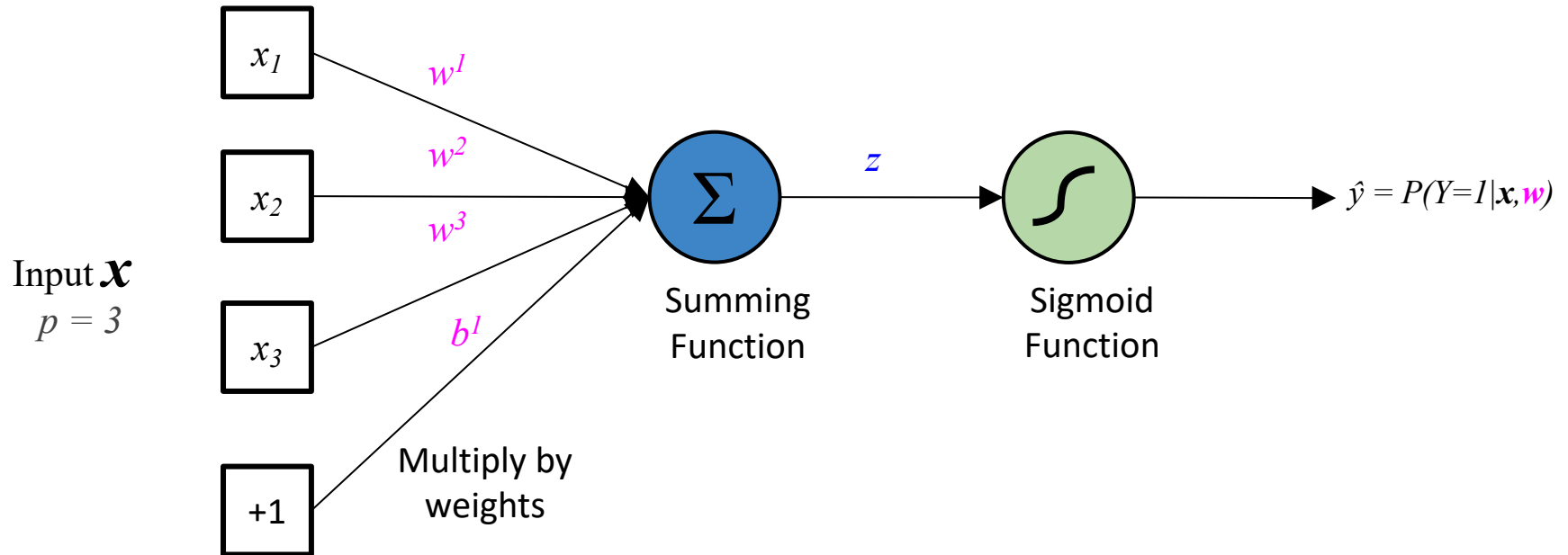# ReWrite Logistic Regression as two stages:

**First**:
Summing

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

**Second**:
Sigmoid
Squashing

$$\hat{y} = P(y=1|x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p}} = \frac{e^z}{1 + e^z}$$

# One "Neuron": Block View of Logistic Regression

Input $\boldsymbol{x}$
$p = 3$

$x_1$

$x_2$

$x_3$

$+1$

$w^1$

$w^2$

$w^3$

$b^1$

Multiply by weights

$\Sigma$

Summing Function

$z$

$\int$

Sigmoid Function

$\hat{y} = P(Y{=}1|\boldsymbol{x},\boldsymbol{w})$

$$z = \boldsymbol{w^T}{\cdot}\boldsymbol{x} + b$$

$$y = sigmoid(z)$$

$$= \frac{e^z}{1 + e^z}$$

46

# e.g., "Block View" of Logistic Regression

$W$ is a vector $\quad\quad\quad z$ is a vector

$x$

Input

$W$

$*$

Dot Product

$z$

Sigmoid

output $\hat{y}$

$E\ (\hat{y}, y)$

loss

parameterized block, W needs to be learned

No Parameters to Learn

# Review: Stochastic GD ➡

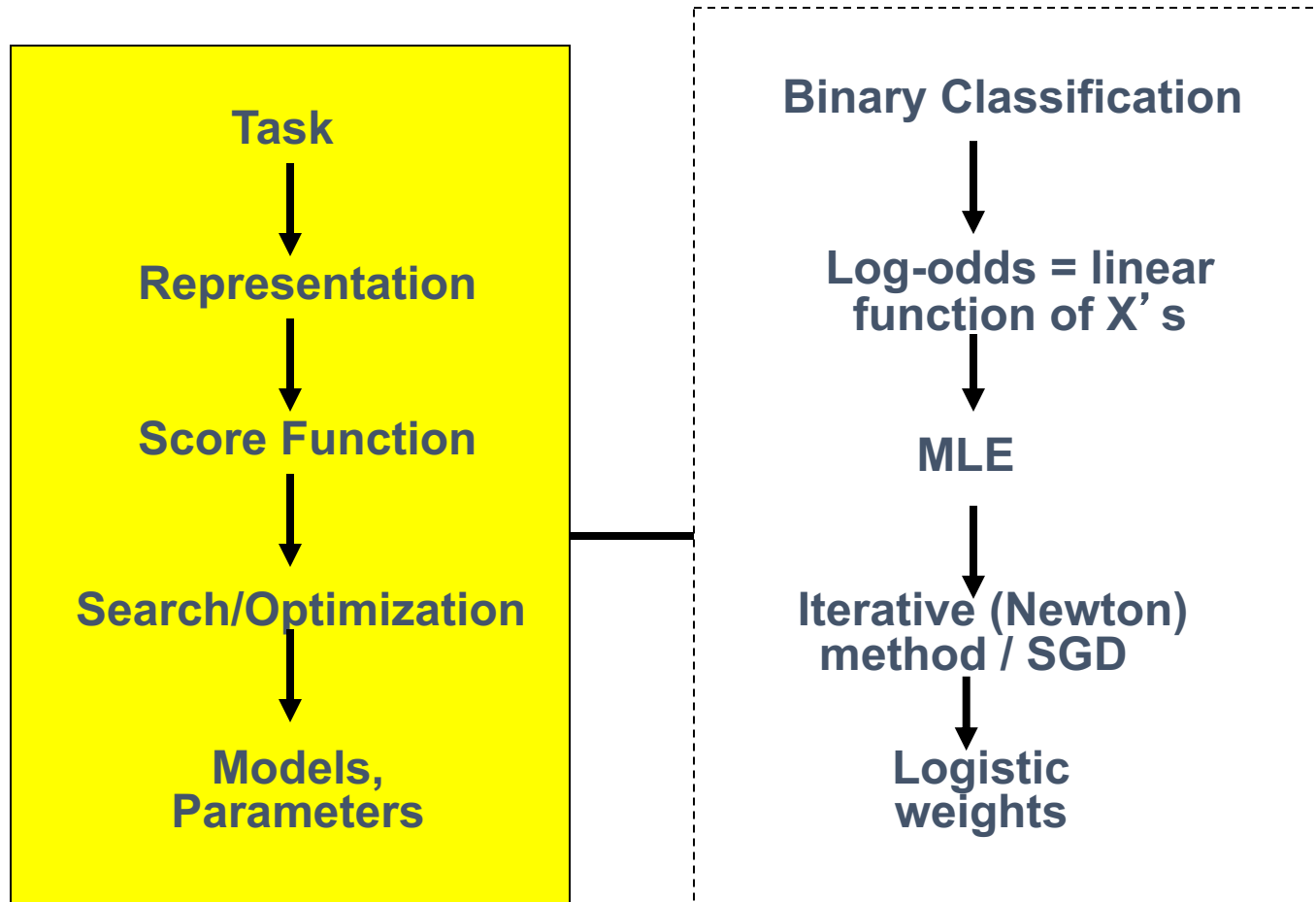- For LR: linear regression, We have the following gradient descent rule:

$$\theta_j^{\,t+1} = \theta_j^{\,t} - \alpha \frac{\partial}{\partial \theta_j} J(\theta)\bigg|_t$$

- ➡ For neural network, we have the delta rule

$$\Delta \mathbf{w} = -\eta \frac{\partial E}{\partial W^t}$$

$$W^{t+1} = W^t - \eta \frac{\partial E}{\partial W^t} = W^t + \Delta w$$

# Logistic Regression



**Task** → **Representation** → **Score Function** → **Search/Optimization** → **Models, Parameters**

**Binary Classification** → **Log-odds = linear function of X's** → **MLE** → **Iterative (Newton) method / SGD** → **Logistic weights**

$$P(y = 1 \mid x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

# Three major sections for classification

• We can divide the large variety of classification approaches into roughly three major types

1. Discriminative

       directly estimate a decision rule/boundary

       e.g~~., support vector machine~~, decision tree, ~~logistic regression,~~

       e.g. neural networks (NN), deep NN

2. Generative:

       build a generative statistical model

       e.g., Bayesian networks, Naïve Bayes classifier

3. Instance based classifiers

       - Use observation directly (no models)

       ~~- e.g. K nearest neighbors~~

# References

❑ Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide

❑ Prof. Andrew Moore's slides

❑ Prof. Eric Xing's slides

❑ Prof. Ke Chen NB slides

❑ Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.