

Note:

- HW#5 submission closed now.
- All homework are completed now 😊

CSC872: PAMI – Kazunori Okada (C) 2025

1

1

Note:

- Continue working on your lit. survey project.
- Presentation format: www.pechakucha.com.
More info later.
- Project report due **in three weeks**
 - Read the assignment thoroughly
 - Late policy will apply.
- Project presentation **in three weeks**
 - Submit your slides two days prior (**5/11, 5pm**) by **email**
 - **Read the assignments again to refresh your mem.**
 - **More details for presentation will be emailed soon**
 - **Presentation in alphabetical order**

CSC872: PAMI – Kazunori Okada (C) 2025

2

2

Classification

Supervised ML

Regression & Learning

CSC 872
Pattern Analysis and Machine Intelligence

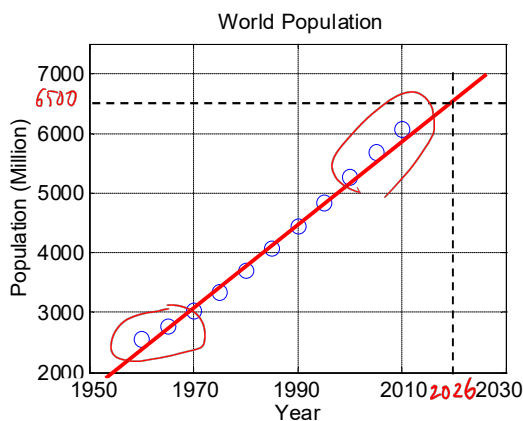
CSC872: PAMI – Kazunori Okada (C) 2025

3

3

What is Regression ?

Statistics 101



What is the population in year 2026

1. Fit a line

2. Find our prediction

- **Regression** is a statistical analysis to find a function representing **input-output relation** from data samples

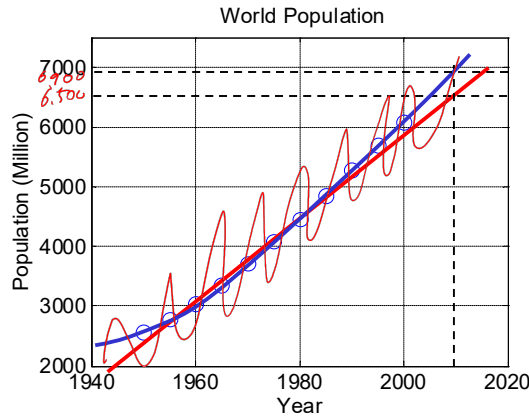
CSC872: PAMI – Kazunori Okada (C) 2025

http://www.populationmedia.org/issues/popgrowth_data.html

4

4

What is Regression ?



What is the population in year 2026?

1. Fit a **curve**

2. Find our prediction

- We must choose **appropriate function form** to do regression

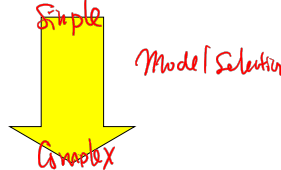
CSC872: PAMI – Kazunori Okada (C) 2025

http://www.populationmedia.org/issues/popgrowth_data.html

5

In Nutshell

- Goal: estimate input-output relation from data
- For: prediction/forecasting/modeling
- You need to
 - 1) *Pick a form of parametric function*
 - Line:
 - Polynomial Curves:
 - General Curves:
 - General Functions:
 - 2) *Fit the function to the data*
 - Maximum likelihood estimation (MLE) is the foundation



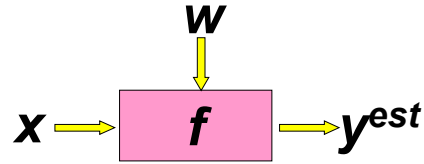
CSC872: PAMI – Kazunori Okada (C) 2025

6

6

Learning Machine Interpretation

- Learning Machine: $y = f(x)$



- Input:** independent variable

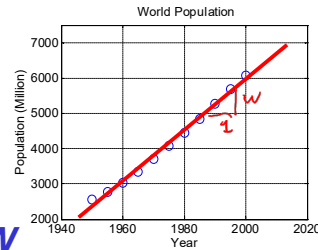
– E.g., $x = \text{year}$

- Output:** dependent variable

– E.g., $y = \text{population}$

- Function:** parameterized by W

– E.g., $f(x,w) = wx$: line



CSC872: PAMI – Kazunori Okada (C) 2025

7

7

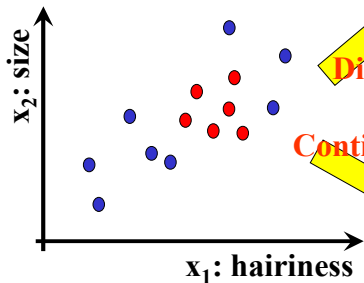
Training Data X (Supervised ML)

- Supervised learning

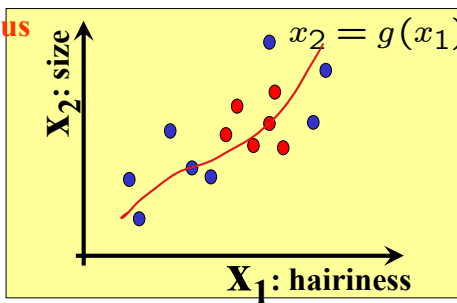
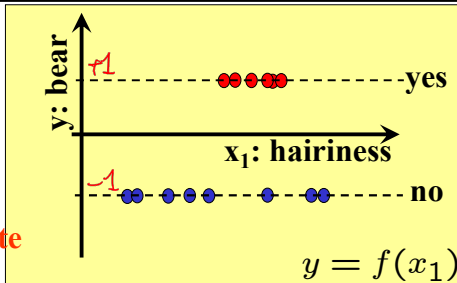
$$y_i = f(x_i, w)$$

$x_i \in \mathbb{R}^M$
 $y_i \in \mathbb{R}$

$$X = \{(x_i, y_i)\}_{i=1}^{N=14}$$



$$y_i = \begin{cases} +1 & \text{if bear} \quad \bullet \\ -1 & \text{otherwise} \quad \bullet \end{cases}$$



CSC872: PAMI – Kazunori Okada (C) 2025

8

8

Regression Types

Depends on the form of parameterized functions

- Linear Regression (line/plane/hyperplane)
- Polynomial Regression (polynomial curve)
- Non-linear Regression (general curve)
- Radial-Basis Function Regression (basis sum)
- Piecewise Linear Regression (line segments)
- Non-parametric Regression (KDE)
- Robust Regression (robust estimation!!!)

CSC872: PAMI – Kazunori Okada (C) 2025

9

9

Regression Types

Depends on the form of parameterized functions

- Linear Regression (line/plane/hyperplane)
- Polynomial Regression (polynomial curve)
- Non-linear Regression (general curve)
- Radial-Basis Function Regression (basis sum)
- Piecewise Linear Regression (line segments)
- Non-parametric Regression (KDE)
- Robust Regression (robust estimation!!!)

CSC872: PAMI – Kazunori Okada (C) 2025

10

10

Linear Regression

- Simplest one parameter case

$$y \in \mathbb{R}, x \in \mathbb{R}, w \in \mathbb{R}$$

- Data is formed by: $y = wx + \text{noise}$

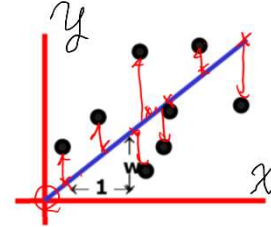
- Unknown scalar w
- Noise is independent **random variable**
- Noise is **normally-distributed** with zero mean & σ^2

$$\hookrightarrow \text{noise model} \longrightarrow \mathcal{N}(0, \sigma^2)$$

- Output y is then also a random variable

$$\text{with } P(y|x, w) = \text{Normal}(\text{mean "wx", variance "\sigma^2"})$$

- Given data: N *i. i. d.* evidences $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
- Regression Problem: *Find w from data such that*



CSC872: PAMI – Kazunori Okada (C) 2025

11

11

Bayesian Linear Regression

- *Find w from data such that it maximizes the **posterior distribution**:*

$$P(w | (x_1, y_1), (x_2, y_2), \dots, (x_N, y_N))$$

- Infer $P(w|\text{data})$ from data likelihood $P(y|x, w)$ using Bayes rules ! $P(w)$
 - Conjugate prior etc, A bit complicated so...

CSC872: PAMI – Kazunori Okada (C) 2025

12

12

Maximum-Likelihood Estimate

- Find w from data such that it maximizes the **data likelihood function**:

$$P(y | x_1, x_2, \dots, x_N, w) = N(y; wx, \sigma^2)$$

- As usual, Let's do some algebra to simplify

Algebra Joy: You know this by now

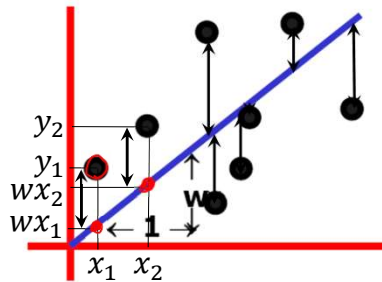
- For what w is this data most likely to have happened?
- For what w , is $P(y_1, \dots, y_N | x_1, \dots, x_N, w)$ maximized?
- For what w , is $\prod_{i=1}^N P(y_i | x_i, w)$ maximized? *↪ i.i.d.*
- For what w , is $\prod_{i=1}^N \exp(-\frac{1}{2} \frac{(y_i - wx_i)^2}{\sigma^2})$ maximized? *↪ $N(wx, \sigma^2)$*
- For what w , is $\sum_{i=1}^N -\frac{1}{2} \frac{(y_i - wx_i)^2}{\sigma^2}$ maximized? *↪ log*
- For what w , is $\sum_{i=1}^N (y_i - wx_i)^2$ minimized? *↪ algebra*
This is known as Least Squares method

Least Squares Method

$$y = wx + \text{noise} \sim \mathcal{N}(0, \sigma^2)$$

- MLE of w is one that minimizes the **sum-of-squares of residuals (errors)**

$$E(w) = \sum_i (y_i - wx_i)^2$$



$$\Sigma \left\{ \begin{array}{l} (y_1 - wx_1)^2 \\ (y_2 - wx_2)^2 \\ \vdots \\ (y_i - wx_i)^2 \\ \vdots \\ (y_N - wx_N)^2 \end{array} \right\}$$

CSC872: PAMI – Kazunori Okada (C) 2025

15

15

Really a quadratic optimization

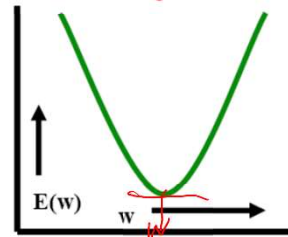
$$E(w) = aw^2 - bw + c$$

- MLE of w is one that minimizes the **sum-of-squares of residuals (errors)**

$$\begin{aligned} E(w) &= \sum_i (y_i - wx_i)^2 = \sum_i (x_i^2 w^2 - 2x_i y_i w + y_i^2) \\ &= \underbrace{\left(\sum_i x_i^2 \right)}_a w^2 - \left(2 \sum_i x_i y_i \right)_b w + \left(\sum_i y_i^2 \right)_c \end{aligned}$$

- We want to minimize a quadratic function of w

$$\frac{\partial E(w)}{\partial w} = 0$$



CSC872: PAMI – Kazunori Okada (C) 2025

16

16

STOP: We have a closed-form solution!

- For linear regression w/ normal-distributed noise
- **MLE = Least Squares !!!**

$$w^{MLE} \equiv \operatorname{argmin}_w E(w) = \operatorname{argmin}_w \sum_i (y_i - wx_i)^2$$

$$\Leftrightarrow \frac{\partial E(w)}{\partial w} = 0$$

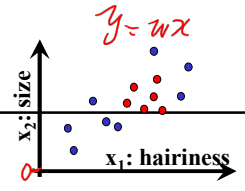
$$w = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

r: Pearson correlation coefficient
 σ_x : standard deviation of $\{x_i\}$
 σ_y : standard deviation of $\{y_i\}$

$$y = wx = \frac{\sum_i x_i y_i}{\sum_i x_i^2} x = \frac{\sigma_y}{\sigma_x} r x$$

Multivariate Case?

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iM} \end{pmatrix}$$



- What if input x is a vector $(x_1, \dots, x_M)^T$?

- Model is $y = w^t x + \epsilon$
 $= w_1 x_1 + w_2 x_2 + \dots + w_m x_m + \dots + w_M x_M + \epsilon$

- Given data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

$$X = \begin{bmatrix} x_1^t \\ x_2^t \\ \vdots \\ x_N^t \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$Y = XW + \epsilon \quad W = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{pmatrix}$$

a set of linear Eq.

- MLE of w is given by

$$w^{MLE} = (X^t X)^{-1} X^t Y \quad (\text{Pseudo Inverse})$$

$$w = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

$$y = w^t x = ((X^t X)^{-1} X^t Y)^t x$$

$$Y = WX$$

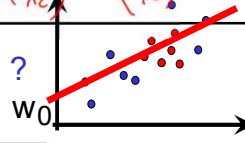
$$X^{-1} Y = X^{-1} W X = W$$

$$\Rightarrow W = X^{-1} Y$$

Constant Term?

Handwritten: $w = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$, $x^t = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix}$

- What if the line does not intersect the origin?



| X_1 | X_2 | Y |
|-------|-------|-----|
| 2 | 4 | 16 |
| 3 | 4 | 17 |
| 5 | 5 | 20 |

Handwritten: $w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$
 $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

| X_0 | X_1 | X_2 | Y |
|-------|-------|-------|-----|
| 1 | 2 | 4 | 16 |
| 1 | 3 | 4 | 17 |
| 1 | 5 | 5 | 20 |

Handwritten: $w^t x + w_0 + \epsilon$

$Y = w_1 X_1 + w_2 X_2 + \epsilon$
 $= w^t x + \epsilon$

$Y = w_0 + w_1 X_1 + w_2 X_2 + \epsilon$
 $= w_0 X_0 + w_1 X_1 + w_2 X_2 + \epsilon$

- Model is

$y = w^t z + \epsilon$,

$z = (1, x^t)^t$

Handwritten: $= w^t \begin{pmatrix} 1 \\ x \end{pmatrix} + \epsilon$

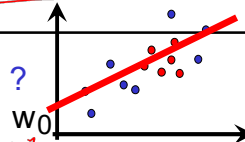
Handwritten: $w_0 = w_0 \times 1$

$= w_0 \cdot 1 + w_1 x_1 + w_2 x_2 + \dots + w_M x_M + \epsilon$

Constant Term?

Handwritten: Linear SVM

- What if the line does not intersect the origin?
- Model is $y = w^t z + \epsilon$, $z = (1, x^t)^t$



- Given data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Handwritten: Add new column of 1s

$z = \begin{bmatrix} z_1^t \\ z_2^t \\ \vdots \\ z_N^t \end{bmatrix} = \begin{bmatrix} (1, x_1^t)^t \\ (1, x_2^t)^t \\ \vdots \\ (1, x_N^t)^t \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1M} \\ 1 & x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{NM} \end{bmatrix}$ $Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

Handwritten: $Y = Z w + \epsilon$

- MLE of w is given by

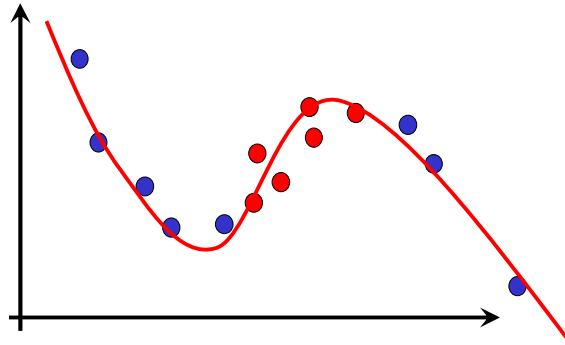
$w = (Z^t Z)^{-1} Z^t Y$

Most general linear regression formula !!

$y = w^t z = ((Z^t Z)^{-1} Z^t Y)^t \begin{bmatrix} 1 \\ x \end{bmatrix}$

Polynomial Regression

- What if I want to fit a polynomial curve?



- You can reuse the linear formula!!!

Quadratic Regression

$$y = wx + \epsilon$$

$$y = w^T x + w_0 + \epsilon$$

$$= w_0 + w_1 x_1 + w_2 x_2 + \epsilon$$

- Model for 2D is

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 + \epsilon$$

$$= \mathbf{w}^t \mathbf{z} + \epsilon, \quad \mathbf{z} = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$$

- Given data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, $\mathbf{x}_i = (x_{i1}, x_{i2})$

$$\mathbf{Z} = \begin{bmatrix} z_1^t \\ z_2^t \\ \vdots \\ z_N^t \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{11}^2 & x_{11}x_{12} & x_{12}^2 \\ 1 & x_{21} & x_{22} & x_{21}^2 & x_{21}x_{22} & x_{22}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N2} & x_{N1}^2 & x_{N1}x_{N2} & x_{N2}^2 \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

You change here

- MLE of \mathbf{w} is given by

$$\mathbf{w} = (\mathbf{Z}^t \mathbf{Z})^{-1} \mathbf{Z}^t \mathbf{Y}$$

But all these are same

$$y = \mathbf{w}^t \mathbf{z} = ((\mathbf{Z}^t \mathbf{Z})^{-1} \mathbf{Z}^t \mathbf{Y})^t \mathbf{z}$$

Qth degree Polynomial Regression

- Model is the same but with different \mathbf{z} *→ long*
 $y = \mathbf{w}^t \mathbf{z} + \epsilon, \quad \mathbf{z} = (1, q(x))$
- Given data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- $q(\mathbf{x})$: all products of powers of inputs up to Qth degree

$$\mathbf{z} = \begin{bmatrix} z_1^t \\ z_2^t \\ \vdots \\ z_N^t \end{bmatrix} = \begin{bmatrix} 1 & q(x_1) \\ 1 & q(x_2) \\ \vdots & \vdots \\ 1 & q(x_N) \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

You cannot model complicated curves

Because Z gets large very quick by increasing the complexity of the curve by increasing Q when more than 1D input

- MLE of \mathbf{w} is given by
 $\mathbf{w} = (\mathbf{Z}^t \mathbf{Z})^{-1} \mathbf{Z}^t \mathbf{Y}$

$$y = \mathbf{w}^t \mathbf{z} = ((\mathbf{Z}^t \mathbf{Z})^{-1} \mathbf{Z}^t \mathbf{Y})^t \mathbf{z}$$

CSC872: PAMI – Kazunori Okada (C) 2025

23

23

Radial Basis Function Regression

- Can we generalize the idea of the Polynomial Regression?
 - Basically, you construct Z with different z with various products of inputs
 - Then, use the same pseudo inverse formula
- Let's construct \mathbf{z} with some function $\phi(\mathbf{x})$ of input \mathbf{x}
- Model is $y = w_0 + w_1 \phi_1(\mathbf{x}) + \dots + w_K \phi_K(\mathbf{x}) + \epsilon$ *Meron Koval*
 $= \mathbf{w}^t \mathbf{z} + \epsilon, \quad \mathbf{z} = (1, B(\mathbf{x}) = \phi_1, \dots, \phi_K)$
- $B(\mathbf{x})$ is called a **basis** whose linear combination gives an output
- We choose the basis to be symmetric about a center \mathbf{c} with spread W then call it **radial basis function**

$$\phi_k(x) = \text{RadialBasisFunction}\left(\frac{|x - c_k|}{W_k}\right)$$
- RBF Regression performs the linear regression with $B(\mathbf{x})$ defined with the radial basis function

CSC872: PAMI – Kazunori Okada (C) 2025

24

24

Non-linear Regression

- What if I want to fit more general nonlinear function $f(\mathbf{x}; \mathbf{w})$?
- Let's do the same as before !
 - Assume a general model of $y = f(x; w) + \epsilon$
 - Normally-distributed independent noise
 - Likelihood $P(\mathbf{Y}|\mathbf{X}, \mathbf{w})$ is $Normal(\text{mean } f(\mathbf{x}; \mathbf{w}), \text{variance } \sigma^2)$
 - MLE of \mathbf{w} = LS of \mathbf{w}

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$$

$$\Leftrightarrow \frac{\partial}{\partial \mathbf{w}} \sum_i (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 = 0$$

$$\Leftrightarrow \sum_i (y_i - f(\mathbf{x}_i; \mathbf{w})) \frac{\partial f(\mathbf{x}_i; \mathbf{w})}{\partial \mathbf{w}} = 0$$

Ooops! how to solve this about \mathbf{w} ???

We are doomed. we are stuck here?...

CSC872: PAMI – Kazunori Okada (C) 2025

25

25

Energy Optimization

- Recall the non-parametric modeling lecture...
- The savior is to go “**iterative**” to solve $w = \operatorname{argmin}_{\mathbf{w}} E(\mathbf{w})$
- Minimizing the **Energy/Error/Cost/Potential function** by
 - Define an **iterative step** $\text{move}(\mathbf{w}, E(\mathbf{w}))$
 - Then find an **initial** solution \mathbf{w}_0
 - Then find a **sequence** $\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_m$ by doing

$$\text{step1: } w_{\text{new}} = w_{\text{old}} + \text{move}(w_{\text{old}}, E(w))$$

$$\text{step2: } w_{\text{old}} = w_{\text{new}}$$

$$\text{step3: } \text{go to step1}$$

- To do this right, you need to find $\text{move}(\mathbf{w}, E(\mathbf{w}))$ so that \mathbf{w}_m converges to a **local minimum** of $E(\mathbf{w})$

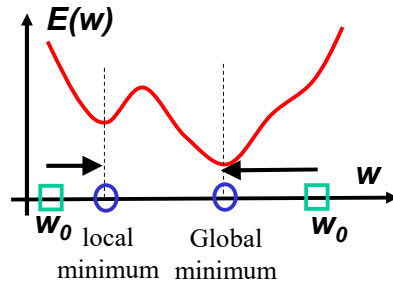
CSC872: PAMI – Kazunori Okada (C) 2025

26

26

In general

- You want $w = \operatorname{argmin}_w E(w)$



Iterative solution can get stuck at local minimum depending on initialization

- What if I have a maximization problem ?

$$w = \operatorname{argmax}_w E(w)$$



$$w = \operatorname{argmin}_w -E(w)$$

You can always interchange minimization and maximization problem

CSC872: PAMI – Kazunori Okada (C) 2025

27

How are we going to solve it?

- Various ways

- Line search

- Hill-Climbing

- – Gradient-Descent (Steepest-Descent/Ascent)

- Conjugate-Gradient

- Levenberg-Marquart

- Newton's Method

- – Simulated Annealing

- EM-algorithm

- Mean Shift

- More and more...

Optimization
We study this during the lectures for search methods in AI

statistical modeling

CSC872: PAMI – Kazunori Okada (C) 2025

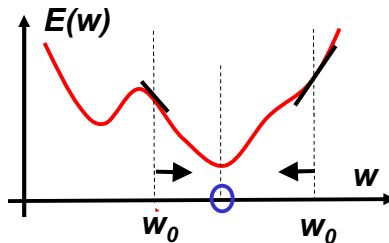
28

28

Gradient-Descent

- We define the step move function as a **negative of partial derivatives** of the energy w.r.t. the unknown parameter W

$$w_{new} = w_{old} - \eta \left. \frac{\partial E(w)}{\partial w} \right|_{w=w_{old}} \quad \eta > 0$$



η is a learning rate set to a small constant (e.g., 0.05)

CSC872: PAMI – Kazunori Okada (C) 2025

29

29

Steepest-Descent/Ascent

- Umm, sorry but, I want to maximize really, want to go up
- DON'T LIKE GOING DOWN!!!**
- Just flip the sign!!!**

$$w_{new} = w_{old} + \eta \left. \frac{\partial E(w)}{\partial w} \right|_{w=w_{old}}$$

Gradient-Descent

- Steepest-Ascent:** a type of greedy iterative search algorithm we learned in our lecture on search for AI
- Umm: the same
- Gradient-Ascent = Steepest-Ascent**

CSC872: PAMI – Kazunori Okada (C) 2025

30

30

Simulated Annealing



- Are there a way to avoid getting stuck in a local minimum?
- **Yes:** called “**simulated annealing**” making it **stochastic**
 - step1: make a **random** move $W_{old} \xrightarrow{\text{random}} W_{new} \quad E \uparrow$
 - step2: take this move if reduces the energy
 - step3: else take it with certain acceptance probability
 - step4: go to step1

Acceptance Prob. $P = \exp\left(\frac{E(W_{old}) - E(W_{new})}{T_i}\right) e^{-\frac{\Delta E}{T}}$

sampling from Maxwell–Boltzman distribution

- When temperature T_i is $\frac{\Delta E}{T}$
 - High: random walk
 - Low: stochastic steepest descent
- Can converge to the global minimum when scheduling a gradual decreasing of the temperature (**cooling schedule**)



31

Summary

- Regression & Learning
 - What is regression?
 - Maximum Likelihood Estimate & Least Squares Method
 - Linear regression
 - Polynomial regression
 - Radial Basis Function regression
 - Gradient-descent
 - Simulated Annealing
- Next
 - Artificial Neural Network
 - End of the LDA FP

32

Artificial Neural Network

- Today's lecture on ML-based regression can be directly used to understand and implement ANN!
- The note will be uploaded in the course web
- Read the note!
 - Neural Network
 - Perceptron
 - Multi-Layered Network
 - Backpropagation
 - Other Networks