# EBU7240
# Computer Vision

## - Introduction to Deep Learning -

*Semester 1, 2021*

**Changjae Oh**

# Outline

- **Machine learning basics++**

- **Introduction to deep learning**

- **Linear classifier**

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# What is Machine Learning?

- Learning = Looking for a Function

Speech Recognition

$$f\left( \text{ } \right) = \text{ "Hello"}$$

Handwritten Recognition

$$f\left( \text{2} \right) = \text{ "2"}$$

Weather forecast

$$f\left( \text{Thursday} \right) = \text{ "Saturday"}$$

Play video games

$$f\left( \text{ } \right) = \text{ "move left"}$$

# Machine Learning: Basic Concept

- **Prediction task**

  - Regression: returns a specific value

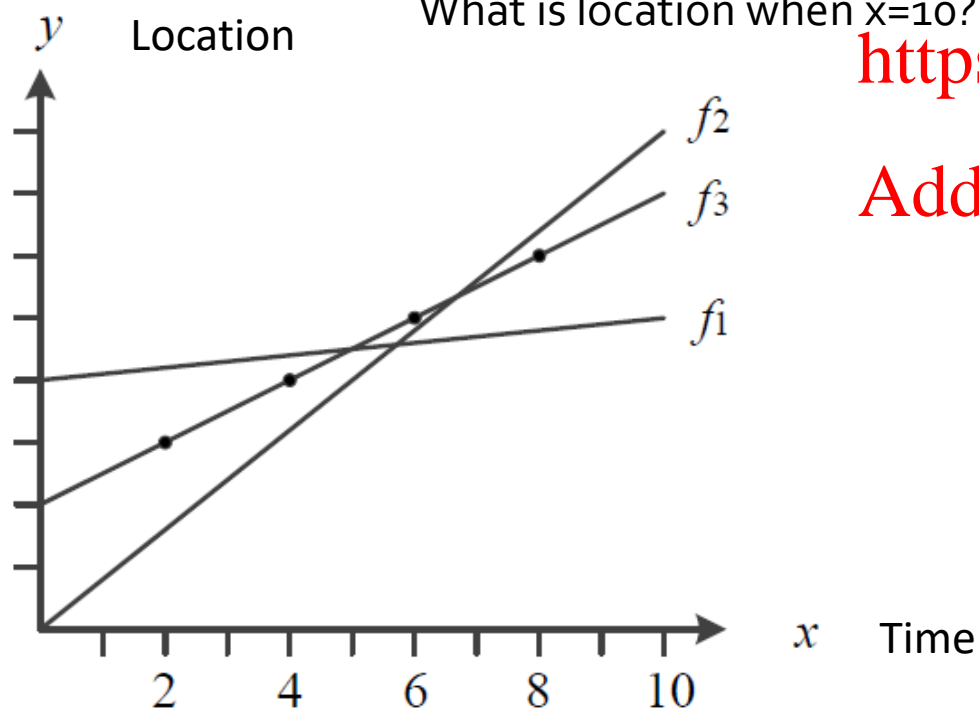  - Classification: returns a class label
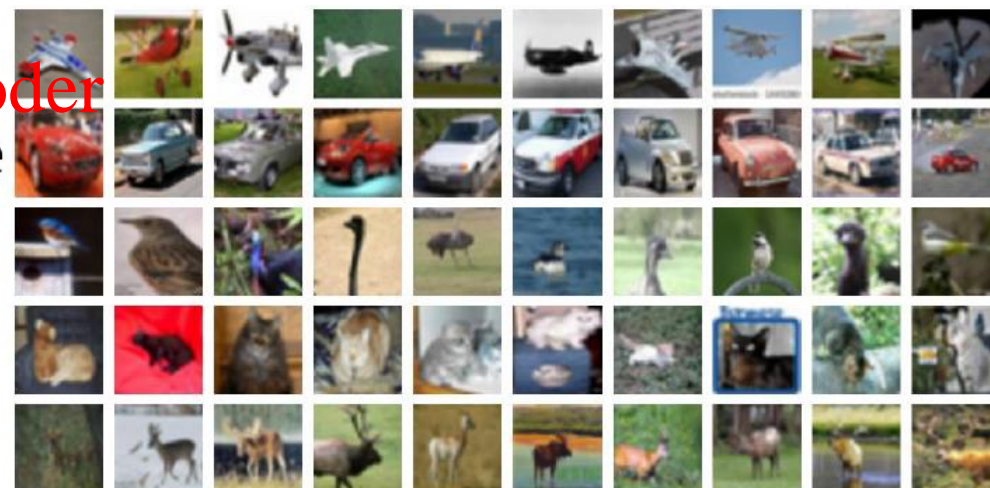
**Regression example**

What is location when x=10?

**Classification example**

For a new image, what class does it belong to?



y  Location

$f2$

$f3$

$f1$

x  Time

2   4   6   8   10



airplane

automobile

bird

cat

deer

# Machine Learning: Basic Concept

- ## Training data

$$\mathbb{X} = \{\mathbf{x}_1 = (2.0), \mathbf{x}_2 = (4.0), \mathbf{x}_3 = (6.0), \mathbf{x}_4 = (8.0)\}$$
$$\mathbb{Y} = \{y_1 = 3.0, y_2 = 4.0, y_3 = 5.0, y_4 = 6.0\}$$

## Training the model with data

- Finding optimal parameters
- Starting at a random value, increasing accuracy to compute optimal parameters

Assignment Project Exam Help

$$y = wx + b$$

https://powcoder.com

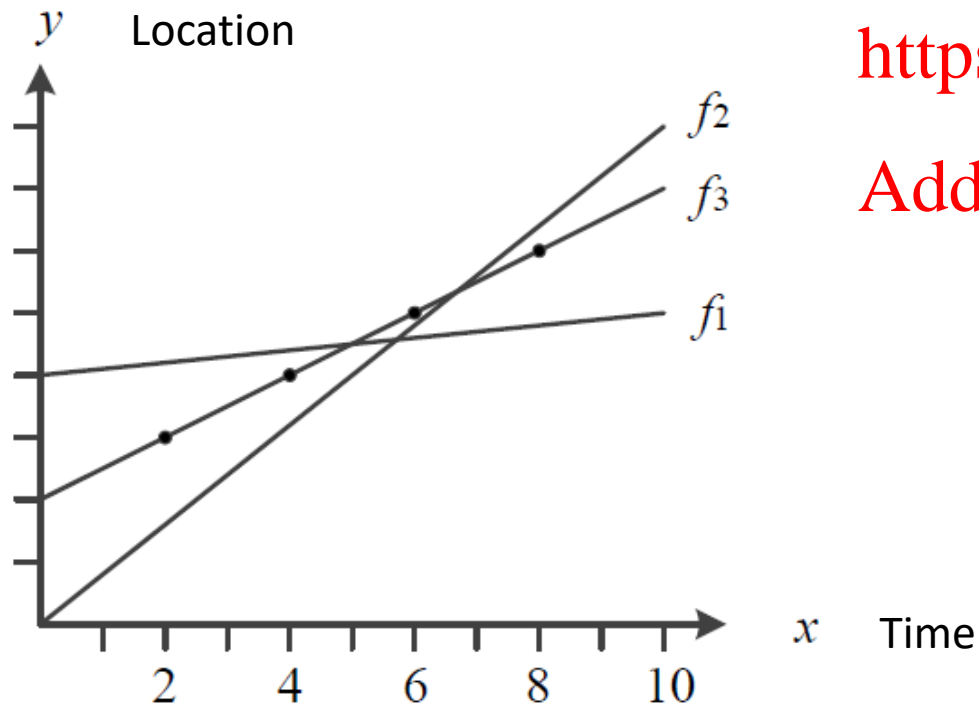*Optimal parameter w=0.5 b=2.0*

Add WeChat powcoder

## Goal

- Minimize errors for new samples (test set)
- Generalization
  refers to high performance for test sets

# Machine Learning: Basic Concept

- **Multi-dimensional feature space**
  - d-dimensional data: $\mathbf{x}=(x_1,x_2, \dots ,x_d)^{\mathrm{T}}$      Note) d=784 in MNIST data

- **Linear classifier for $d$-dimensional data**

  - 1-D linear classifier          # of variables $= d + 1$
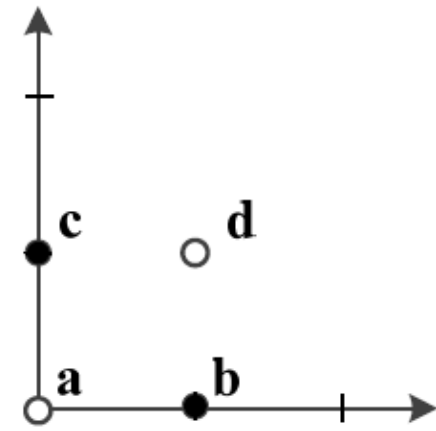    - Widely used in machine learning

  $$y = w_1 x_1 + w_2 x_2 + \cdots + w_d x_d + b$$

  - 2-D linear classifier          # of variables $= \frac{(d+1)(d+2)}{2}$

$$y = w_1 x_1^2 + w_2 x_2^2 + \cdots + w_d x_d^2 + w_{d+1} x_1 x_2 + \cdots + w_{\frac{d(d+1)}{2}} x_{d-1} x_d + w_{\frac{d(d+1)}{2}+1} x_1 \cdots + w_{\frac{d(d+1)}{2}+d} x_d + b$$

# Machine Learning: Basic Concept

- **Feature space transformation**
  - Map a linearly non-separable feature space into separable space

**Toy example**

$$\mathbf{x} = (x_1, x_2)^T \longrightarrow \mathbf{x}' = \left( \frac{x_1}{2x_1 x_2 + 0.5}, \frac{x_2}{2x_1 x_2 + 0.5} \right)^T$$

$\mathbf{a} = (0,0)^T \longrightarrow \mathbf{a}' = (0,0)^T$

$\mathbf{b} = (1,0)^T \longrightarrow \mathbf{b}' = (2,0)^T$

$\mathbf{c} = (0,1)^T \longrightarrow \mathbf{c}' = (0,2)^T$

$\mathbf{d} = (1,1)^T \longrightarrow \mathbf{d}' = (0.4,0.4)^T$

Original feature space

Transformed feature space, which is linearly separable

# Machine Learning: Basic Concept
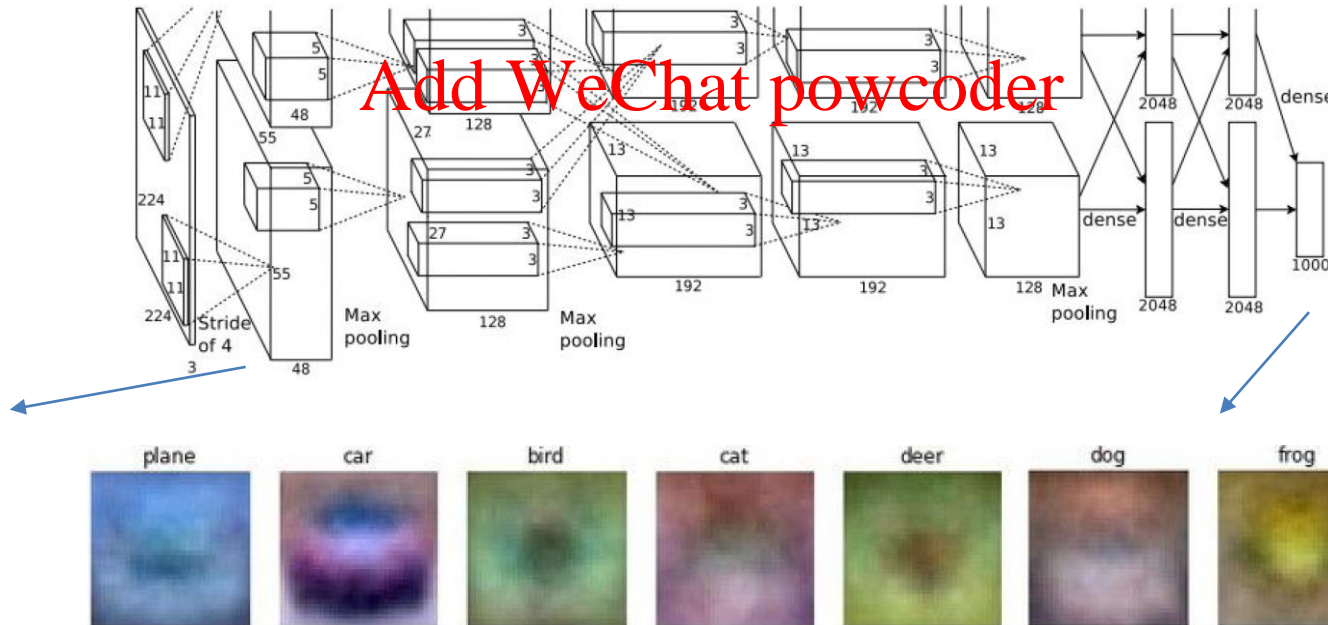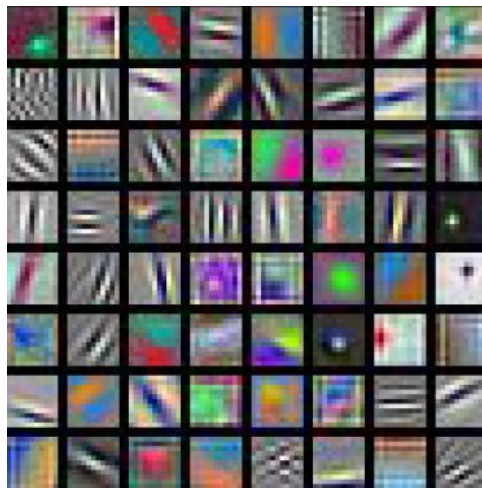
- **Representation learning**
  - Aims to find good feature space automatically
  - Deep learning finds a hierarchical feature space by using neural networks with multiple hidden layers.
    - The first hidden layer has low-level features (edge, corner points, etc.), and the right-hand side features advanced features (face, wheel, etc.)

# Data for Machine Learning

- **The quality of the training data**

  - To increase estimation accuracy, diverse and enough data should be collected for a given application.

  - Ex) After learning from a database with a frontal face only, the recognition accuracy of side face will be degraded.
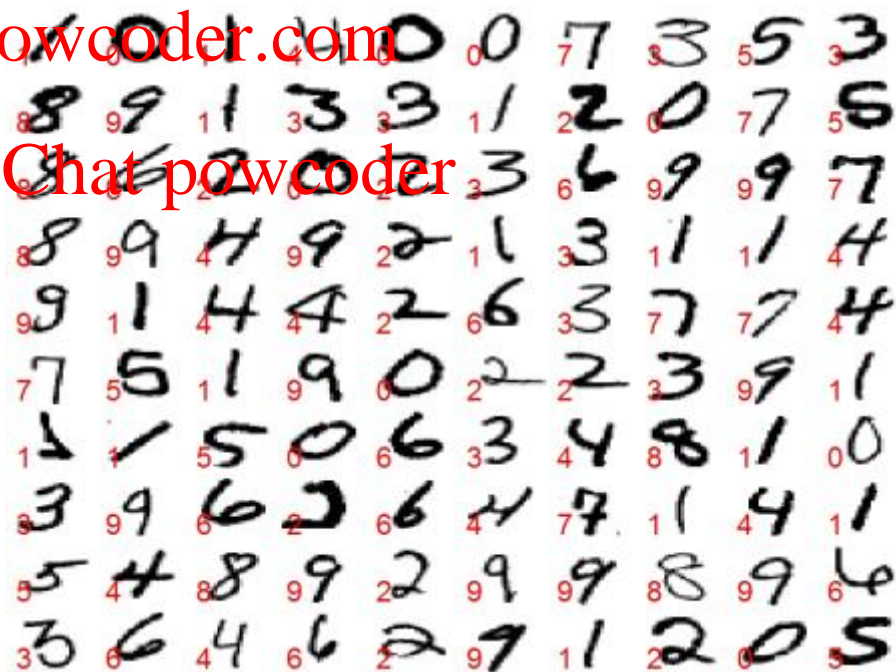
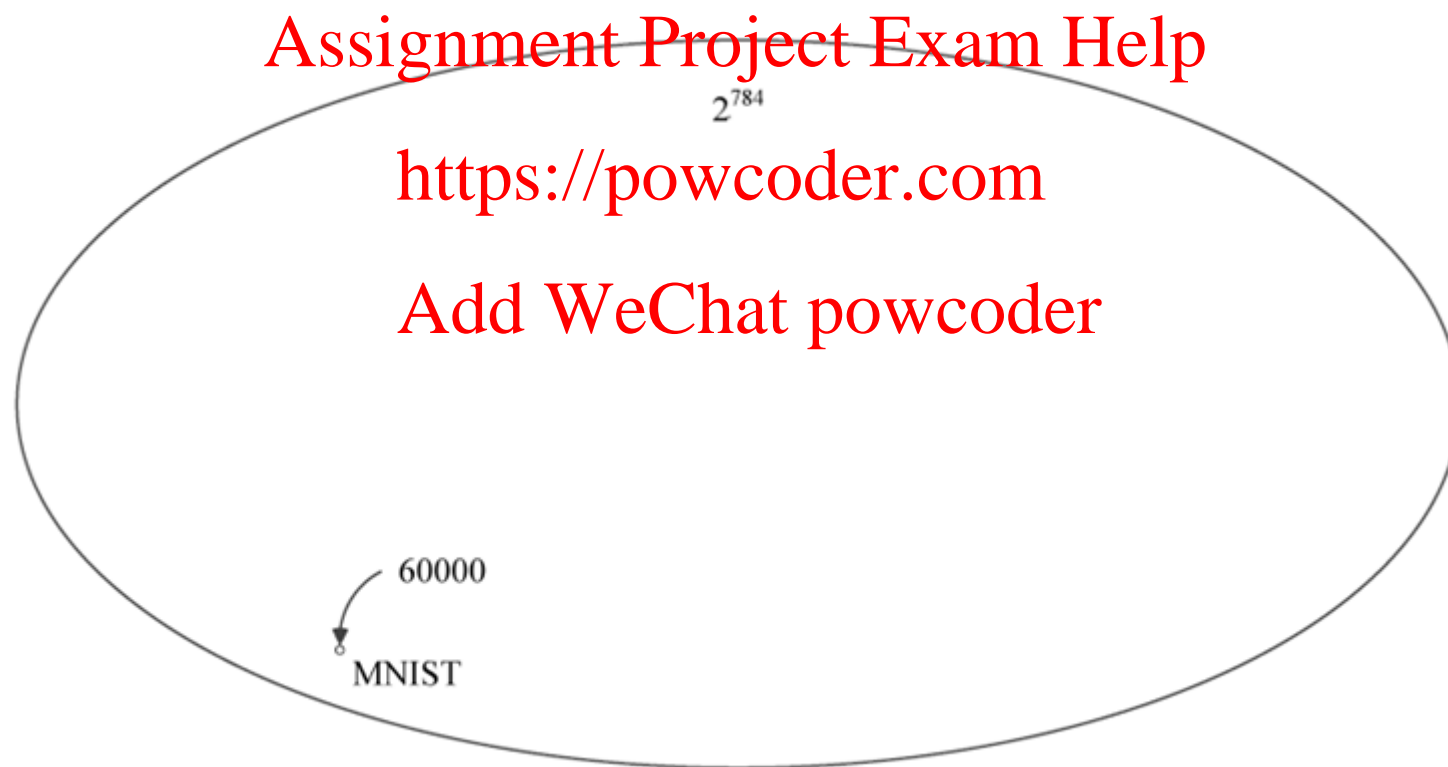- **MNIST database**

  - Handwritten numeric database

  - Training data: 60,000

  - Test data: 10,000

# Data for Machine Learning

- **Database size vs. training accuracy**
  - Ex) MNIST: 28*28 binary image
    → The total number of possible samples is $2^{784}$, but MNIST has 60,000 training images.

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

$2^{784}$

60000

MNIST

# Data for Machine Learning

- **How does a small database achieve high performance?**

  - In a feature space, the actual data is generated in a very small subspace

     ,     is unlikely to happen.

    Assignment Project Exam Help

  - Manifold assumption

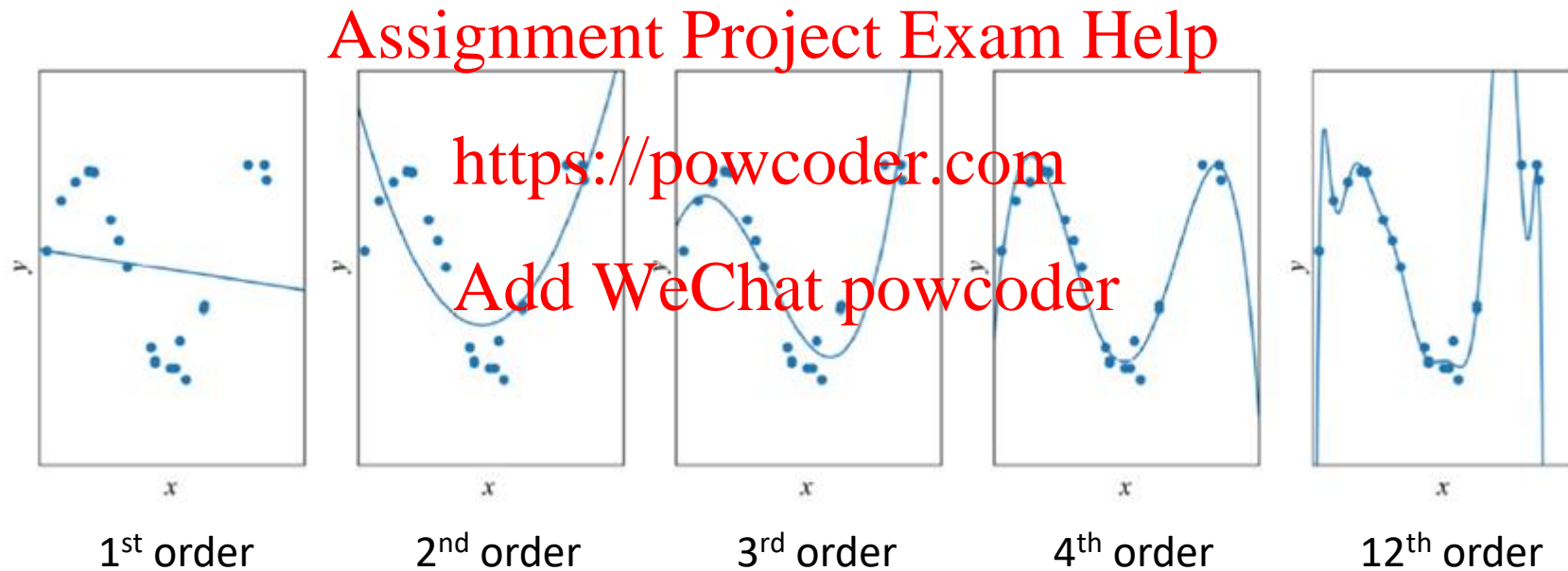    https://powcoder.com

    Smooth change according to certain rules like 

    Add WeChat powcoder

# Training Model: under-fitting vs. over-fitting

- **Under-fitting**
    - Model capacity is too small to fit the data accordingly.
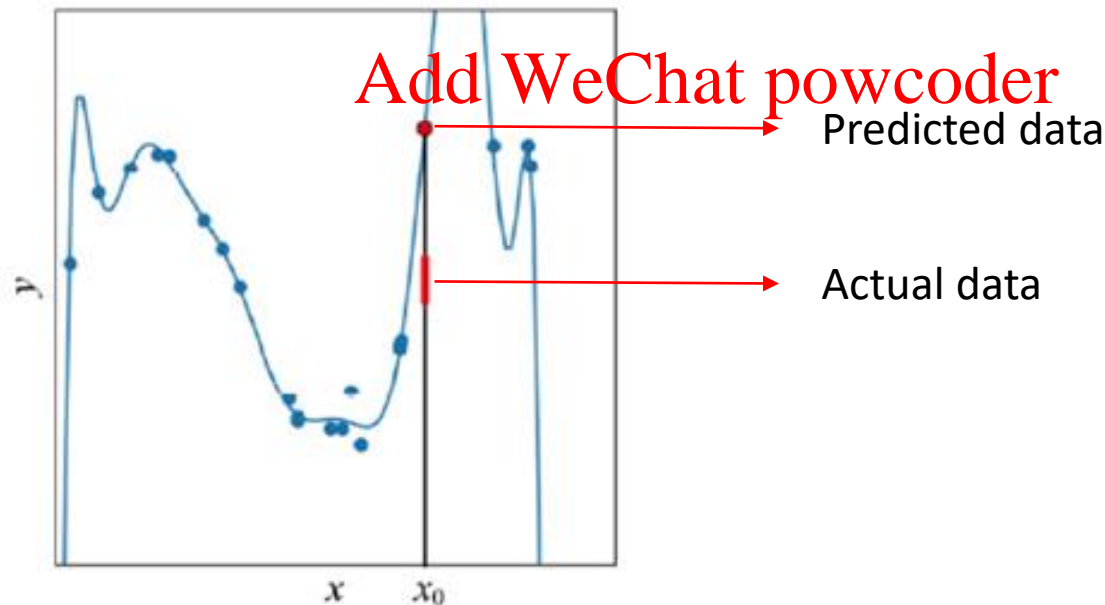    - Model with higher order can be used.

| 1st order | 2nd order | 3rd order | 4th order | 12th order |

# Training Model: under-fitting vs. over-fitting

- **Over-fitting**
  - $12^{th}$ order polynomial model approximates perfectly for the training set.
  - But if you anticipate "new" data, there's a big problem.
  - Since the model capacity is large, the training process also accepts data noise.
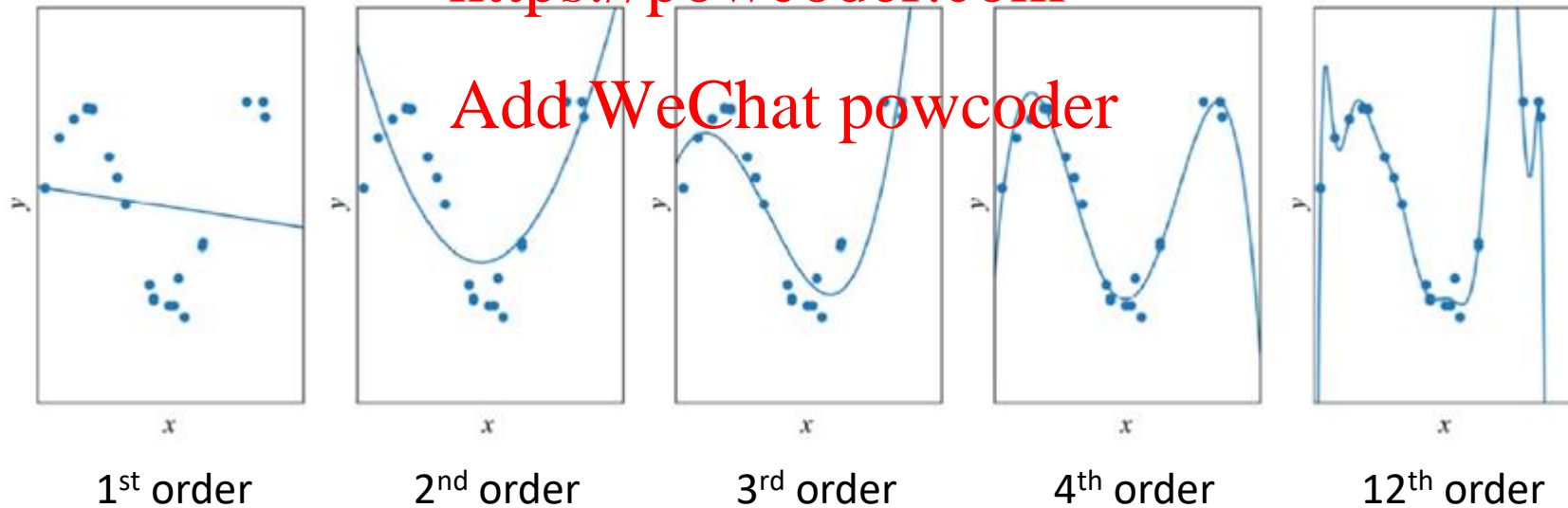  - The model with the appropriate capacity should be selected.



Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Predicted data

Actual data

# Training Model: under-fitting vs. over-fitting

- $1^{st}$ and $2^{nd}$ order model show poor performance for both the training and the test set.

- $12^{th}$ order model shows high performance in training set, but low performance in test set. $\rightarrow$ low generalization ability

- $3^{rd}$ and $4^{th}$ order model are lower than the $12^{th}$ order model for the training set, but the test set has high performance. $\rightarrow$ higher generalization capability

| 1st order | 2nd order | 3rd order | 4th order | 12th order |

# Spectrum of supervision



Computer vision

Supervised learning

Semi-Supervised learning

Unsupervised learning

Reinforcement learning

# Spectrum of supervision

- **Supervised learning**
  - Both the feature vector $\mathbb{X}$ and the output $\mathbb{Y}$ are given.
  - Regression and classification problem

- **Unsupervised learning**
  - The feature vector $\mathbb{X}$ is given, but the output $\mathbb{Y}$ is not given.
  - Ex) Clustering, density estimation, feature space conversion

# Spectrum of supervision

- **Reinforcement learning**
  - The output is given, but it is different from supervised learning.
  - Ex) Go
    - Once the game is over, you get a point (credit).
      If you win, get 1, and -1 otherwise.
    - The credit should be distributed to each sample of the game.

- **Semi-supervised Learning**
  - Some of data have both $\mathbb{X}$ and $\mathbb{Y}$, but others have only $\mathbb{X}$.
  - It is becoming important, since it is easy to collect $\mathbb{X}$, but $\mathbb{Y}$ requires manual tasks.

# Introduction to Deep Learning

# From Wiki

- **Deep learning**
  - is a branch of machine learning based on a set of algorithms that attempt to model high-level abstractions in data by using multiple processing layers, with complex structures or otherwise, composed of multiple non-linear transformations.

# Deep learning success

- **Image classification**
- Machine translation
- Speech recognition
- Speech synthesis
- Game playing
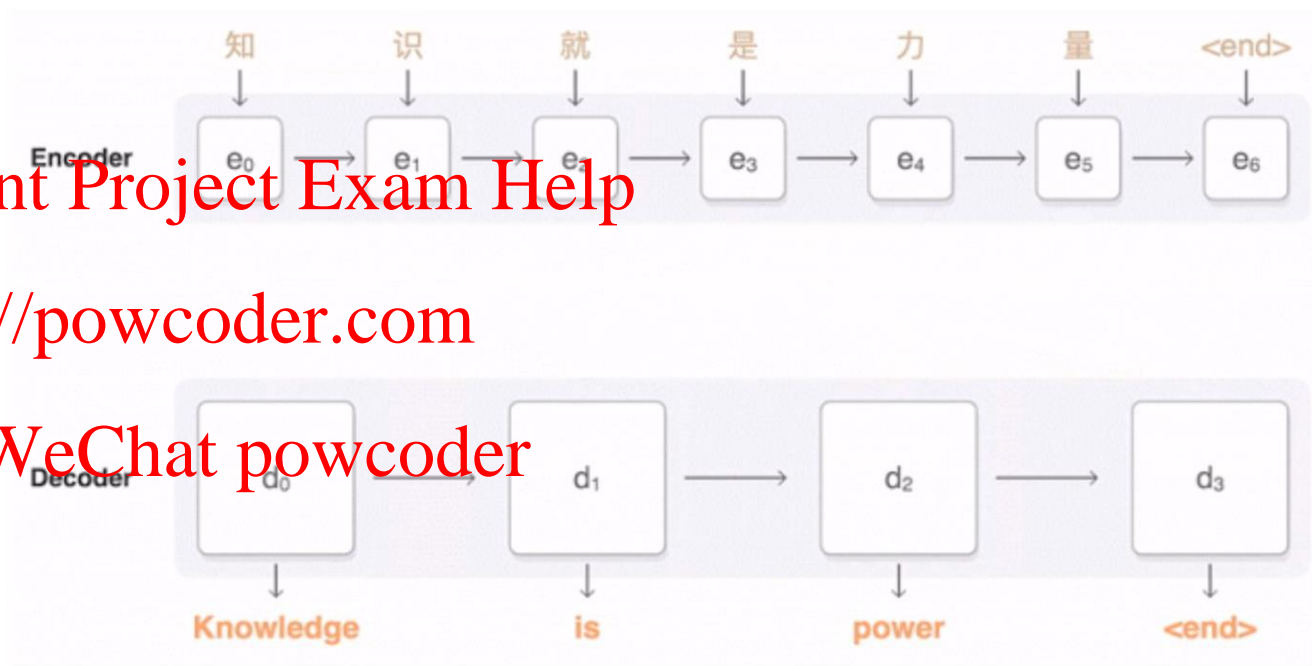
- .. and many, many more

# Deep learning success

- Image classification
- **Machine translation**
- Speech recognition
- Speech synthesis
- Game playing

- .. and many, many more

# Deep learning success

- Image classification
- Machine translation
- **Speech recognition**
- Speech synthesis
- Game playing

- .. and many, many more
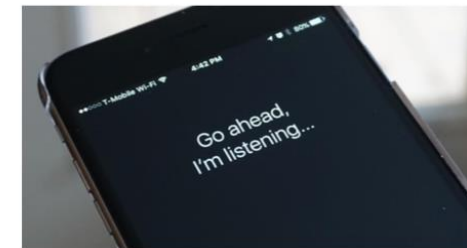
Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Deep learning success

- Image classification
- Machine translation
- Speech recognition
- **Game playing**



- .. and many, many more

# Why deep learning?

- **Hand-crafted features vs. Learned features**



Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Why now?

- Large datasets

- GPU hardware advances + Price decreases

- Improved techniques (algorithm)

# What is Deep Learning?

- **Stacked Functions Learned by Machine**
  - **End-to-end training**: what each function should do is learned automatically
  - Deep learning usually refers to neural network based model



Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# What is Deep Learning?

- **Stacked Functions Learned by Machine**
  - **Representation Learning:** learning features/representations
  - **Deep Learning:** learning (multi-level)  features and an output



*Features / Representations*

# Machine Learning vs. Deep Learning

- **Deep vs Shallow: Image Recognition**
  - Shallow model using machine learning

# Machine Learning vs. Deep Learning

- **Deep vs Shallow: Image Recognition**
  - Deep model using deep learning



Features / Representations

# Machine Learning vs. Deep Learning

- **Machine Learning vs. Deep Learning**

**Machine Learning = Feature descriptor + Classifier**

Feature: <u>hand-crafted</u> domain-specific knowledge
Describing your data with features that computer can understand
Ex) SIFT, Bag-of-Words (BoW), Histogram of Oriented Gradient (HOG)

Optimizing the classifier weights on features
Ex) Nearest Neighbor (NN), Support Vector Machine (SVM), Random Forest (RF)

# Machine Learning vs. Deep Learning

- **Machine Learning vs. Deep Learning**

**Deep Learning** = Feature descriptor + Classifier

Feature: Representation learned by machine
Automatically learned internal knowledge

Optimizing the classifier weights on features

⇨ **Neural network based model**

A series of linear classifiers and non-linear activations + Loss function

# Deep Learning

- **A single neuron**



$$z = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b = \sum_{k=1}^{d} w_k x_k + b \qquad y = \frac{1}{1 + e^{-z}}$$

# Deep Learning

- A single layer with multiple neurons



$$z_1 = \boldsymbol{w}_1^{\mathrm{T}}\boldsymbol{x} + b_1 = \sum_{k=1}^{d} w_{1k}x_k + b_1 \qquad y_1 = \frac{1}{1 + e^{-z_1}}$$

$$z_M = \boldsymbol{w}_M^{\mathrm{T}}\boldsymbol{x} + b_M = \sum_{k=1}^{d} w_{Mk}x_k + b_M \qquad y_M = \frac{1}{1 + e^{-z_n}}$$

# Deep Learning

- **Deep Neural Network**
  - Cascading the neurons to form a neural network



Each layer consists of the linear classifier and activation function

# Linear classifier

Neural Network

Linear classifiers

# Parametric Approach

- **(Review) Unit 3 ML basics and classification**

Image

**3072X1**

$$f(x,W) = Wx + b$$

Assignment Project Exam Help

**10X1**    **10X3072**

https://powcoder.com

$f(\mathbf{x},\mathbf{W})$

Add WeChat powcoder

**10X1**

**10** numbers giving class scores

Array of **32x32x3** numbers
(3072 numbers total)

**W**

parameters
(or weights)

# Parametric Approach

- **(Review) Unit 3 recognition**

  Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Stretch pixels into column

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

W

| 56 |
|----|
| 231 |
| 24 |
| 2 |

**+**

| 1.1 |
|-----|
| 3.2 |
| -1.2 |

**=**

| -96.8 | Cat score |
|-------|-----------|
| 437.9 | Dog score |
| 61.95 | Ship score |

Input image

# What do we need now?

- **Functions to measuring the error between the output of a classifier and the given target value.**

  - Let's talk about designing error (a.k.a. loss) functions!

# EBU7240
# Computer Vision
## - Loss Functions -

*Semester 1, 2021*

**Changjae Oh**

# Loss function

- **Loss function**
  - quantifies our unhappiness with the scores across the training data.

- **Type of loss function**
  - Hinge loss
  - Cross-entropy loss
  - Log likelihood loss
  - Regression loss

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Loss Function: Hinge Loss

- **Binary hinge loss (=binary SVM loss)**

$$L_i = \max(0, 1 - y_i \cdot s)$$

$$s = \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x_i} + b$$

$y_i = \pm 1$ for positive/negative samples

- **Hinge loss (=multiclass SVM loss)**
  - $C$: The number of class $(> 2)$

$$L_i = \sum_{j=1, j \neq y_i}^{C} \max(0, s_j - s_{y_i} + 1)$$

$\boldsymbol{x_i}$: input data (e.g. image)
$y_i$: class label (integer, $1 \leq y_i \leq C$)

$$\boldsymbol{s} = \mathbf{W}\boldsymbol{x_i} + \boldsymbol{b}$$

$$\mathbf{W} = \begin{pmatrix} \boldsymbol{w}_1^T \\ \boldsymbol{w}_2^T \\ \vdots \\ \boldsymbol{w}_C^T \end{pmatrix} \qquad \boldsymbol{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_C \end{pmatrix}$$

# Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, \mathbf{W}) = \mathbf{W}x + b$ are

Given a dataset of examples
$$\{(x_i, y_i)\}_{i=1}^N$$

$x_i$: image

$y_i$: class label (integer)

Loss over the dataset is a
sum of loss over examples:

| | | | |
|------|------|------|------|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

# Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx + b$ are

**Multiclass SVM loss (=hinge loss)**

$$L_i = \sum_{i \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

where score vector $s = f(x_i, W)$

|      |       |      |      |
|------|-------|------|------|
| cat  | **3.2** | 1.3  | 2.2  |
| car  | 5.1   | **4.9** | 2.5  |
| frog | -1.7  | 2.0  | **-3.1** |

# Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx + b$ are

**Multiclass SVM loss (=hinge loss)**

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 5.1 - 3.2 + 1) + max(0, -1.7 - 3.2 + 1)
= max(0, 2.9) + max(0, -3.9)
= 2.9 + 0 = 2.9

= max(0, 1.3 - 4.9 + 1) + max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0 = 0

= max(0, 2.2 - (-3.1) + 1) + max(0, 2.5 - (-3.1) + 1)
= max(0, 6.3) + max(0, 6.6)
= 6.3 + 6.6

|      |      |      |      |
|------|------|------|------|
| cat  | **3.2** | 1.3 | 2.2 |
| car  | 5.1  | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Loss | 2.9  | 0    | 12.9 |

Loss over full dataset is average

**L = (2.9 + 0 + 12.9)/3 = 5.27**

# Loss Function: Log Likelihood Loss

- ## Log likelihood loss

$$p = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_C \end{pmatrix}$$ probability for $i^{th}$ image
(It is assumed to be *normalized*, i.e. $|p| = 1$.)

$$L_i = -\log p_j \text{ where } j \text{ satisfies } z_{ij} = 1$$

$z_i$: class label for $i^{th}$ image
($C \times 1$ vector, $z_{ij} = 1$ when $j = y_i$ and 0 otherwise)

$y_i$: class label (integer, $1 \leq y_i \leq C$)

**Example**
Suppose $i^{th}$ image belongs to class 2 and $C = 10$.

$$z_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad p = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \implies L_i = -\log 0.7$$

# Loss Function: Cross-entropy Loss

- **Cross-entropy loss**

$$L_i = -\sum_{j=1}^{C}\left(z_{ij}\log p_j + (1-z_{ij})\log(1-p_j)\right)$$

$$p = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_C \end{pmatrix}$$

probability for $i^{th}$ image
(It is assumed to be *normalized*, i.e. $|p| = 1$.)

$z_i$: class label for $i^{th}$ image
($C \times 1$ vector, $z_{ij} = 1$ when $j = y_i$ and 0 otherwise)

$y_i$: class label (integer, $1 \leq y_i \leq C$)

**Example**

Suppose $i^{th}$ image belongs to class 2 and $C = 10$.

$$z_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad p = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \implies L_i = -\log(1-0.1) - \log 0.7 - \log(1-0.2)$$

# Softmax Activation Function

- Softmax activation function

**scores = unnormalized log probabilities of the classes.**

→ **Probability can be computed using scores as below.**

Probability of class label being $k$ for an image $x_i$

$$P(Y = k | X = x_i) = p_k = \boxed{\frac{e^{s_k}}{\sum_{j=1}^{C} e^{s_j}}}$$

<span style="color:red">Softmax activation function</span>

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">https://powcoder.com</span>

unnormalized probabilities <span style="color:red">Add WeChat powcoder</span>

$x_i$: image

$y_i$: class label (integer, $1 \leq y_i \leq C$)

cat **3.2**
car 5.1
frog -1.7

exp →

**24.5**
164.0
0.18

normalize →

**0.13**
0.87
0.00

unnormalized log probabilities

probabilities

$$s = W x_i + b$$

$$W = \begin{pmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_C^T \end{pmatrix}$$

# Softmax + Log Likelihood Loss

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_{j=1}^{C} e^{s_j}}\right)$$

unnormalized probabilities

| | | | |
|---|---|---|---|
| cat | **3.2** | **24.5** | **0.13** |
| car | 5.1 | 164.0 | 0.87 |
| frog | -1.7 | 0.18 | 0.00 |

exp → normalize →

L_i = -log(0.13)
= **0.89**

unnormalized log probabilities

probabilities

Softmax + Log likelihood loss:
is often called 'softmax classifier'

# Softmax + Cross-entropy Loss

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_{j=1}^{C} e^{s_j}}\right) - \sum_{k=1, k \neq y_i}^{C} \log\left(1 - \frac{e^{s_k}}{\sum_{j=1}^{C} e^{s_j}}\right)$$

unnormalized probabilities

| | | | |
|---|---|---|---|
| cat | **3.2** | **24.5** | **0.13** |
| car | 5.1 | 164.0 | 0.87 |
| frog | -1.7 | 0.18 | 0.00 |

exp → normalize →

unnormalized log probabilities

probabilities

L_i = -log(0.13)-log(1-0.87)-log(1-0.0)
   = 0.89*2

# Loss Function: Regression Loss

- **Regression loss**
  - Using L1 or L2 norms
  - Widely used in pixel-level prediction (e.g. image denoising)

$$L_i = |\boldsymbol{y}_i - \boldsymbol{s}_i|$$

$$L_i = (\boldsymbol{y}_i - \boldsymbol{s}_i)^2$$

$$\boldsymbol{y}_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad \boldsymbol{s}_i = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \implies L_i = |\boldsymbol{y}_i - \boldsymbol{s}_i| = |0 - 0.1| + |1 - 0.7| + |0 - 0.2|$$

# Regularization

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \qquad s = Wx + b$$

Suppose that we found a W such that L = 0. Is this W unique?

**No! 2W is also has L = 0!**

| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

0

**Before:**
= max(0, 1.3 - 4.9 + 1) + max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0 = 0

**With W twice as large:**
= max(0, 2.6 - 9.8 + 1) + max(0, 4.0 - 9.8 + 1)
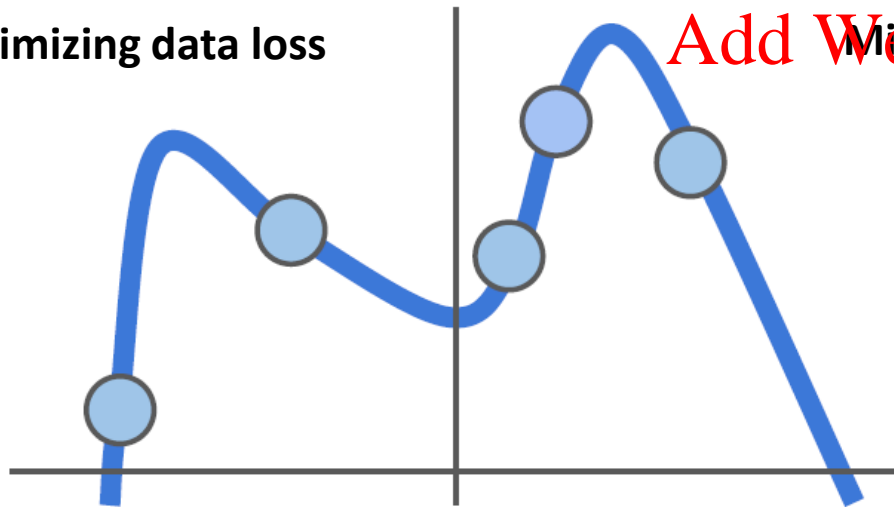= max(0, -6.2) + max(0, -4.8)
= 0 + 0 = 0

# Regularization

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i(f(\boldsymbol{x}_i, \mathbf{W}), y_i) + \lambda R(\mathbf{W})$$

**Data loss**: Model predictions should match training data

**Regularization**: Model should be "simple" to avoid *overfitting*, so it works on test data

**Minimizing data loss**

**Minimizing data + regularization loss**

# Regularization

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i(f(\boldsymbol{x}_i, \mathbf{W}), y_i) + \lambda R(\mathbf{W})$$

$\lambda$: regularization strength (hyperparameter)

- L2 regularization: $R(\mathbf{W}) = \sum_{k,l} W_{k,l}^2$

- L1 regularization: $R(\mathbf{W}) = \sum_{k,l} |W_{k,l}|$

- Elastic net (L1 + L2): $R(\mathbf{W}) = \sum_{k,l} \beta W_{k,l}^2 + |W_{k,l}|$

- Max norm regularization: $|\boldsymbol{w}_j^T| < c$ **for all** $j$

- Dropout (will see later)

- Batch normalization, stochastic depth (will see later)

$$\mathbf{W} = \begin{pmatrix} \boldsymbol{w}_1^T \\ \boldsymbol{w}_2^T \\ \vdots \\ \boldsymbol{w}_C^T \end{pmatrix}$$

# Optimization: Gradient Descent

- **Gradient Descent**
  - The simplest approach to minimizing a loss function

$$\mathbf{W}^{T+1} = \mathbf{W}^T - \alpha \frac{\partial L}{\partial \mathbf{W}^T}$$

Assignment Project Exam Help

  - $\alpha$: step size (a.k.a. learning rate)

https://powcoder.com

Add WeChat powcoder

$$y = f(x) = x^2 - 4x + 3$$
$$y' = f'(x) = 2x - 4$$

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# Optimization: Gradient Descent



original W

negative gradient direction

# Optimization: Stochastic Gradient Descent (SGD)

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i(f(\boldsymbol{x}_i, \mathbf{W}), y_i) + \lambda R(\mathbf{W})$$

Full sum is too expensive when N is large!

$$\frac{\partial L}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial L_i(f(\boldsymbol{x}_i, \mathbf{W}), y_i)}{\partial \mathbf{W}} + \lambda \frac{\partial R(\mathbf{W})}{\partial \mathbf{W}}$$
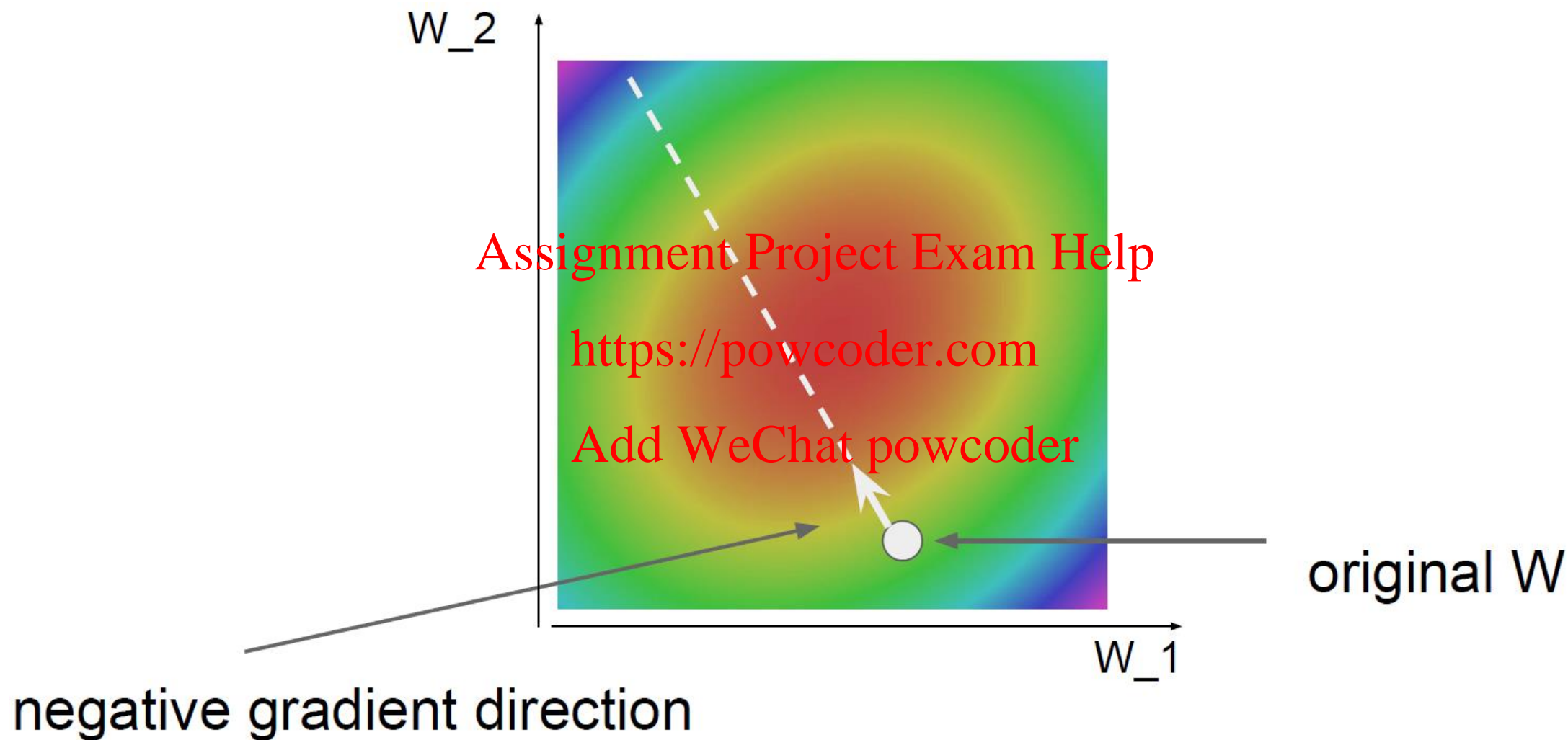
Instead, approximating sum using a **minibatch** of 32 / 64 / 128/ 256 examples is common

```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# EBU7240
# Computer Vision
## - Backpropagation-

*Semester 1, 2021*

**Changjae Oh**

# Backpropagation

- **A widely used algorithm for training feedforward neural networks.**

- **A way of computing gradients of expressions through recursive applicati on of chain rule.**

  – Backpropagation computes the gradient of the loss function with respect to the weig hts of the network (model) for a single input–output example.

- **Gradient Descent**

  – The simplest approach to minimizing a loss function

$$\mathbf{W}^{T+1} = \mathbf{W}^T - \alpha \frac{\partial L}{\partial \mathbf{W}^T}$$

  – $\alpha$: step size (a.k.a. learning rate)

$f''(x)$

$f(x)$

$x$

$f'(x)$

$y = f(x) = x^2 - 4x + 3$

$y' = f'(x) = 2x - 4$

# Derivative

- **Optimization using derivative**

  - 1st order derivative

  $$f'(x) = \lim_{\Delta x \to 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

  - $f'(x)$: The slope of the function, indicating the direction in which the value increases

    → The minima of the objective function may exist in the direction of $-f'(x)$.

    → Gradient descent algorithm: $-f'(x)$

$f''(x)$

$f(x)$

$f'(x)$

$y = f(x) = x^2 - 4x + 3$

$y' = f'(x) = 2x - 4$

$$\mathbf{W}^{T+1} = \mathbf{W}^T - \alpha \frac{\partial L}{\partial \mathbf{W}^T}$$

# Derivative

- **Partial derivative**
    - Derivatives of functions with multiple variables
    - Gradient: the vector of the partial derivative

Ex) $\nabla f, \dfrac{\partial f}{\partial \mathbf{x}}, \left( \dfrac{\partial f}{\partial x_1}, \dfrac{\partial f}{\partial x_2} \right)^{\mathrm{T}}$

$$f(\mathbf{x}) = f(x_1, x_2) = \left( 4 - 2.1x_1^2 + \frac{x_1^4}{3} \right) x_1^2 + x_1 x_2 + (-4 + 4x_2^2)x_2^2$$

$$\nabla f = f'(\mathbf{x}) = \frac{\partial f}{\partial \mathbf{x}} = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)^{\mathrm{T}} = (2x_1^5 - 8.4x_1^3 + 8x_1 + x_2, \; 16x_2^3 - 8x_2 + x_1)^{\mathrm{T}}$$

# Derivative

- ## Jacobian matrix

  - 1$^{st}$ order partial derivative matrix for $\mathbf{f}: \mathbb{R}^d \mapsto \mathbb{R}^m$

$$J = \begin{pmatrix} \dfrac{\partial f_1}{\partial x_1} & \dfrac{\partial f_1}{\partial x_2} & \cdots & \dfrac{\partial f_1}{\partial x_d} \\ \dfrac{\partial f_2}{\partial x_1} & \dfrac{\partial f_2}{\partial x_2} & \cdots & \dfrac{\partial f_2}{\partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial f_m}{\partial x_1} & \dfrac{\partial f_m}{\partial x_2} & \cdots & \dfrac{\partial f_m}{\partial x_d} \end{pmatrix}$$

Ex) $\mathbf{f}: \mathbb{R}^2 \mapsto \mathbb{R}^3 \quad \mathbf{f(x)} = (2x_1 + x_2^2, -x_1^2 + 3x_2, 4x_1 x_2)^T$

$$J = \begin{pmatrix} 2 & 2x_2 \\ -2x_1 & 3 \\ 4x_2 & 4x_1 \end{pmatrix} \quad J|_{(2,1)^T} = \begin{pmatrix} 2 & 2 \\ -4 & 3 \\ 4 & 8 \end{pmatrix}$$

- ## Hessian matrix

  - 2$^{nd}$ order partial derivative matrix

$$H = \begin{pmatrix} \dfrac{\partial^2 f}{\partial x_1 x_1} & \dfrac{\partial^2 f}{\partial x_1 x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_1 x_n} \\ \dfrac{\partial^2 f}{\partial x_2 x_1} & \dfrac{\partial^2 f}{\partial x_2 x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial^2 f}{\partial x_n x_1} & \dfrac{\partial^2 f}{\partial x_n x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_n x_n} \end{pmatrix}$$

Ex) $f(\mathbf{x}) = f(x_1, x_2)$

$$= \left(4 - 2.1x_1^2 + \frac{x_1^4}{3}\right) x_1^2 + x_1 x_2 + (-4 + 4x_2^2)x_2^2$$

$$H = \begin{pmatrix} 10x_1^4 - 25.2x_1^2 + 8 & 1 \\ 1 & 48x_2^2 - 8 \end{pmatrix}$$

$$H|_{(0,1)^T} = \begin{pmatrix} 8 & 1 \\ 1 & 40 \end{pmatrix}$$

# Derivative

- **Chain rule**

$$f(x) = g(h(x))$$
$$f(x) = g(h(i(x)))$$

$$f'(x) = g'(h(x))h'(x)$$
$$f'(x) = g'(h(i(x)))h'(i(x))i'(x)$$

Ex) $f(x) = 3(2x^2 - 1)^2 - 2(2x^2 - 1) + 5$    $h(x) = 2x^2 - 1$

$f'(x) = \underbrace{(3 * 2(2x^2 - 1) - 2)}_{g'(h(x))} \underbrace{(2 * 2x)}_{h'(x)}$ $= 48x^3 - 32x$

# Why are we talking about derivatives?

- **Gradient Descent**
  - The simplest approach to minimizing a loss function

$$\mathbf{W}^{T+1} = \mathbf{W}^{T} - \alpha \frac{\partial L}{\partial \mathbf{W}^{T}}$$

  - $\alpha$: step size (a.k.a. learning rate)

$y = f(x) = x^2 - 4x + 3$

$y' = f'(x) = 2x - 4$

```
# Vanilla Gradient Descent

while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

# Derivative

- **Example) Applying chain rule to single-layer perceptron**

  - Example of composite function

  - Back-propagation: use the chain rule to compute $\frac{\partial L}{\partial \mathbf{W}}$ and $\frac{\partial L}{\partial \boldsymbol{b}}$

Ground truth $\boldsymbol{z}$, $n \times 1$



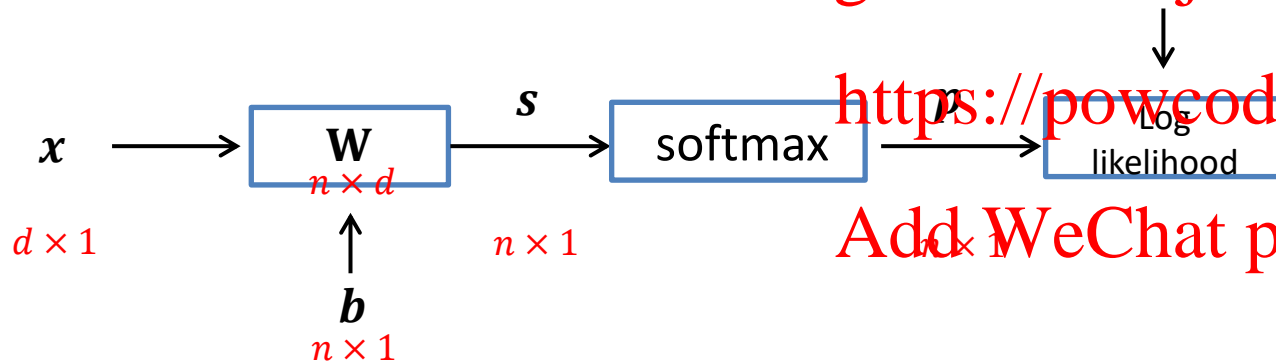$$\boldsymbol{x} \quad \boxed{\begin{array}{c}\mathbf{W} \\ n \times d\end{array}} \quad \xrightarrow{\boldsymbol{s}} \quad \boxed{\text{softmax}} \quad \xrightarrow{\boldsymbol{p}} \quad \boxed{\begin{array}{c}\text{Log} \\ \text{likelihood}\end{array}}$$

$d \times 1$

$n \times 1$

$\uparrow$
$\boldsymbol{b}$
$n \times 1$

$$\frac{\partial L}{\partial \boldsymbol{p}} \implies \frac{\partial L}{\partial \boldsymbol{s}} = \frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}}\frac{\partial L}{\partial \boldsymbol{p}} \implies \frac{\partial L}{\partial \boldsymbol{w}_j} = \frac{\partial \boldsymbol{s}}{\partial \boldsymbol{w}_j}\frac{\partial L}{\partial \boldsymbol{s}} \qquad \frac{\partial L}{\partial \mathbf{W}} = \left(\frac{\partial L}{\partial \boldsymbol{w}_1} \quad \frac{\partial L}{\partial \boldsymbol{w}_2} \quad \cdots \quad \frac{\partial L}{\partial \boldsymbol{w}_n}\right)^{\mathrm{T}}$$

$$\frac{\partial L}{\partial \boldsymbol{b}} = \frac{\partial \boldsymbol{s}}{\partial \boldsymbol{b}}\frac{\partial L}{\partial \boldsymbol{s}} = \frac{\partial L}{\partial \boldsymbol{s}}$$

# Analytic Gradient: Linear Equation

$$s = \mathbf{W}x + b \quad \longleftrightarrow \quad \begin{aligned} s_1 &= w_1^{\mathrm{T}} x + b_1 \\ s_2 &= w_2^{\mathrm{T}} x + b_2 \\ &\vdots \\ s_n &= w_n^{\mathrm{T}} x + b_n \end{aligned}$$

$$s = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} \quad \mathbf{W} = \begin{pmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_n^T \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ & & \vdots & \\ w_{n1} & w_{n2} & \cdots & w_{nd} \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

$$\frac{\partial s_1}{\partial w_1} = x$$

$$\frac{\partial s_2}{\partial w_1} = 0$$

$$\vdots$$

$$\frac{\partial s_n}{\partial w_1} = 0$$

$$\Longrightarrow \quad \frac{\partial s}{\partial w_1} = [x\ 0\ 0\ \cdots 0] \in \Re^{d \times n}$$

j$^{th}$ column

$$\frac{\partial s}{\partial w_j} = [0\ 0\ x\ \cdots 0] \in \Re^{d \times n}$$

$$\frac{\partial s}{\partial b} = \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix} = \mathbf{I} \in \Re^{n \times n}$$

# Analytic Gradient: Linear Equation

$$s = \mathbf{W}x + b \quad \longleftrightarrow \quad \begin{aligned} s_1 &= \boldsymbol{w}_1^{\mathrm{T}}\boldsymbol{x} + b_1 \\ s_2 &= \boldsymbol{w}_2^{\mathrm{T}}\boldsymbol{x} + b_2 \\ &\vdots \\ s_n &= \boldsymbol{w}_n^{\mathrm{T}}\boldsymbol{x} + b_n \end{aligned}$$

$$s = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} \quad \mathbf{W} = \begin{pmatrix} \boldsymbol{w}_1^T \\ \boldsymbol{w}_2^T \\ \vdots \\ \boldsymbol{w}_n^T \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ & & \vdots & \\ w_{n1} & w_{n2} & \cdots & w_{nd} \end{pmatrix} \quad \boldsymbol{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \boldsymbol{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

$$\frac{\partial s_1}{\partial \boldsymbol{x}} = \boldsymbol{w}_1$$

$$\frac{\partial s_2}{\partial \boldsymbol{x}} = \boldsymbol{w}_2 \qquad \Longrightarrow \qquad \frac{\partial \boldsymbol{s}}{\partial \boldsymbol{x}} = [\boldsymbol{w}_1 \ \boldsymbol{w}_2 \ \cdots \boldsymbol{w}_n] = \mathbf{W}^{\mathrm{T}} \in \Re^{d \times n}$$

$$\vdots$$

$$\frac{\partial s_n}{\partial \boldsymbol{x}} = \boldsymbol{w}_n$$

# Analytic Gradient: Sigmoid Function

- **Sigmoid function**

For a scalar $x$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \rightarrow \quad \frac{\partial \sigma(x)}{\partial x} = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x} - 1}{1 + e^{-x}} \frac{1}{1 + e^{-x}} = (1 - \sigma(x))\sigma(x)$$

Similarly, for a vector $\boldsymbol{s} \in \Re^{n \times 1}$

$$\boldsymbol{p} = \sigma(\boldsymbol{s}) = \frac{1}{1 + e^{-\boldsymbol{s}}} \quad \rightarrow \quad \frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = diag\big((1 - \sigma(s_j))\sigma(s_j)\big) = \begin{bmatrix} (1 - \sigma(s_1))\sigma(s_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & (1 - \sigma(s_n))\sigma(s_n) \end{bmatrix}$$

$$\text{for } j = 1, \dots, n$$

# Analytic Gradient: Softmax Activation Function

- **Softmax function**

$$p_k = \frac{e^{s_k}}{\sum_{j=1}^{n} e^{s_j}} \quad \Rightarrow \quad \boldsymbol{p} = \frac{e^{\boldsymbol{s}}}{\sum_{j=1}^{n} e^{s_j}} \quad \text{in vector form}$$

score function $\quad \boldsymbol{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}$

probability $\quad \boldsymbol{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{pmatrix}$

- **1st order derivative of softmax function**

$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = \frac{diag(e^{\boldsymbol{s}}) \cdot \sum e^{s_j} - e^{\boldsymbol{s}}(e^{\boldsymbol{s}})^{\mathrm{T}}}{(\sum e^{s_j})^2} = \frac{1}{(\sum e^{s_j})^2} \left\{ \begin{pmatrix} e^{s_1}\sum e^{s_j} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & e^{s_n}\sum e^{s_j} \end{pmatrix} - \begin{pmatrix} e^{s_1}e^{s_1} & \cdots & e^{s_1}e^{s_n} \\ \vdots & \ddots & \vdots \\ e^{s_n}e^{s_1} & \cdots & e^{s_n}e^{s_n} \end{pmatrix} \right\}$$

# Analytic Gradient: Softmax Activation Function

$$\mathbf{D} = \frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = \frac{diag(e^{\boldsymbol{s}}) \cdot \sum e^{s_j} - e^{\boldsymbol{s}}(e^{\boldsymbol{s}})^{\mathrm{T}}}{(\sum e^{s_j})^2} = \frac{1}{(\sum e^{s_j})^2} \left\{ \begin{pmatrix} e^{s_1}\sum e^{s_j} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & e^{s_n}\sum e^{s_j} \end{pmatrix} - \begin{pmatrix} e^{s_1}e^{s_1} & \cdots & e^{s_1}e^{s_n} \\ \vdots & \ddots & \vdots \\ e^{s_n}e^{s_1} & \cdots & e^{s_n}e^{s_n} \end{pmatrix} \right\}$$

For $a = b$

$$\frac{e^{s_a}(\sum e^{s_j} - e^{s_a})}{(\sum e^{s_j})^2} = p_a(1 - p_a)$$

For $a \neq b$

$$-\frac{e^{s_a}e^{s_b}}{(\sum e^{s_j})^2} = -p_a p_b$$

$$D_{ab} = p_a(\delta_{ab} - p_b)$$

$$\delta_{ab} = \begin{cases} 1 & a = b \\ 0 & \text{otherwise} \end{cases}$$

# Analytic Gradient: Hinge Loss

- **1$^{st}$ order derivative of binary hinge loss**

$$s = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b$$

$$L = \max(0, 1 - y \cdot s)$$
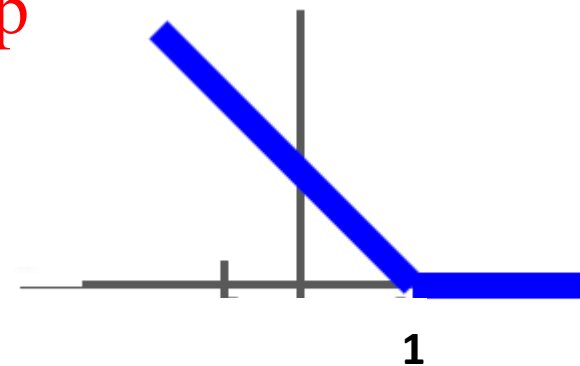
$y = \pm 1$ for positive/negative samples

$$= \begin{cases} 1 - y \cdot s & \text{if } 1 - y \cdot s > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \quad \frac{\partial L}{\partial s} = \begin{cases} -y & \text{if } 1 - y \cdot s > 0 \\ 0 & \text{otherwise} \end{cases}$$

**1**

# Analytic Gradient: Hinge Loss

- **1$^{st}$ order derivative of hinge loss**

$$s = \mathbf{W}x + b$$

$$L = \sum_{j=1, j \neq y}^{n} \max(0, s_j - s_y + 1)$$

Assignment Project Exam Help

$x$: image

$y$: class label (integer, $1 \leq y \leq n$)

https://powcoder.com

Add WeChat powcoder

$$\mathbf{W} = \begin{pmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_n^T \end{pmatrix} \qquad s = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}$$

$$\frac{\partial L}{\partial s_y} = - \sum_{j=1, j \neq y}^{n} 1(s_j - s_y + 1 > 0) \qquad \text{for } j = y$$

$$1(F) = \begin{cases} 1 & \text{if } F \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial L}{\partial s_j} = 1(s_j - s_y + 1 > 0) \qquad \text{for } j \neq y$$

# Analytic Gradient: Log Likelihood Loss

$L = -\log p_y$ where $y$ satisfies $z_y = 1$

$$\boldsymbol{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{pmatrix}$$ probability for $i^{th}$ image
(It is assumed to be *normalized*, i.e. $|\boldsymbol{p}| = 1$.)

$$\frac{\partial L}{\partial \boldsymbol{p}} = -\begin{bmatrix} 0 \\ \vdots \\ 1/p_y \\ \vdots \\ 0 \end{bmatrix}$$ y$^{th}$ row

$y$: class label for $i^{th}$ image ($1 \le y \le n$)

$\boldsymbol{z}$: class probability for $i^{th}$ image

$\boldsymbol{z} = (z_1, z_2 \ldots z_n)^T$, $z_y = 1$ and $z_{k \ne y} = 0$

**Example**

Suppose $i^{th}$ image belongs to class 2 and $n = 10$. $\rightarrow$ $y = 2$

$$\boldsymbol{z} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad \boldsymbol{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \qquad \Rightarrow \qquad \frac{\partial L}{\partial \boldsymbol{p}} = -\begin{bmatrix} 0 \\ 1/0.7 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# Analytic Gradient: Cross-entropy Loss

For simplicity of notation, the index of training image $i$ is omitted here

$$L = -\sum_{j=1}^{n}\left(z_j \log p_j + (1 - z_j)\log(1 - p_j)\right)$$

$$\boldsymbol{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{pmatrix}$$

probability for $i^{th}$ image
(It is assumed to be *normalized*, i.e. $|\boldsymbol{p}| = 1$.)

$$\frac{\partial L}{\partial \boldsymbol{p}} = -\begin{bmatrix} 1/(1 - p_1) \\ \vdots \\ 1/p_y \\ \vdots \\ 1/(1 - p_n) \end{bmatrix}$$

$y^{th}$ row

Assignment Project Exam Help

$y$: class label for $i^{th}$ image ($1 \leq y \leq n$)

$\boldsymbol{z}$: class probability for $i^{th}$ image

https://powcoder.com $z_n)^T$, $z_y = 1$ and $z_{k \neq y} = 0$

Add WeChat powcoder

**Example**

Suppose $i^{th}$ image belongs to class 2 and $n = 10$. $\rightarrow$ $y = 2$

$$\boldsymbol{z} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad \boldsymbol{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \qquad \frac{\partial L}{\partial \boldsymbol{p}} = -\begin{bmatrix} 1/(1 - 0.1) \\ 1/0.7 \\ 0 \\ \vdots \\ 1/(1 - 0.2) \end{bmatrix}$$

# Analytic Gradient: Regression Loss

- **Regression loss**

$$L = (\boldsymbol{y} - \boldsymbol{s})^2 \quad \Rightarrow \quad \frac{\partial L}{\partial \boldsymbol{s}} = -2(\boldsymbol{y} - \boldsymbol{s}) \qquad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \qquad \boldsymbol{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}$$
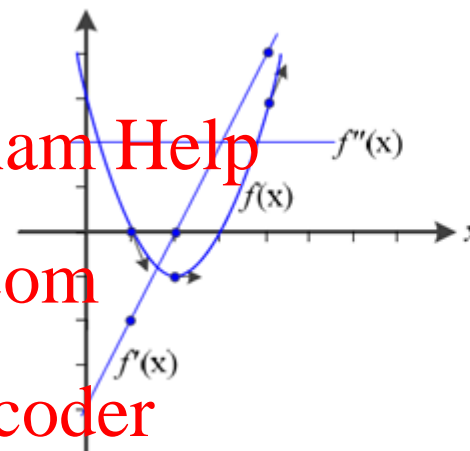
# Why are we talking about derivatives?

- **Gradient Descent**
  - The simplest approach to minimizing a loss function

$$\mathbf{W}^{T+1} = \mathbf{W}^{T} - \alpha \frac{\partial L}{\partial \mathbf{W}^{T}}$$

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

$y = f(x) = x^2 - 4x + 3$

$y' = f'(x) = 2x - 4$

  - $\alpha$: step size (a.k.a. learning rate)

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```