

EBU7240

Computer Vision

- Introduction to Deep Learning -

Semester 1, 2021

Changjae Oh

Outline

- Machine learning basics++
- Introduction to deep learning
- Linear classifier

What is Machine Learning?

- Learning = Looking for a Function

Speech Recognition

$$f(\text{[waveform plot]}) = \text{"Hello"}$$

Handwritten Recognition

$$f(\text{[handwritten digit]}) = \text{"2"}$$

Weather forecast

$$f(\text{[sun icon]} \quad \text{Thursday}) = \text{“} \text{[rain icon]} \text{ Saturday”}$$

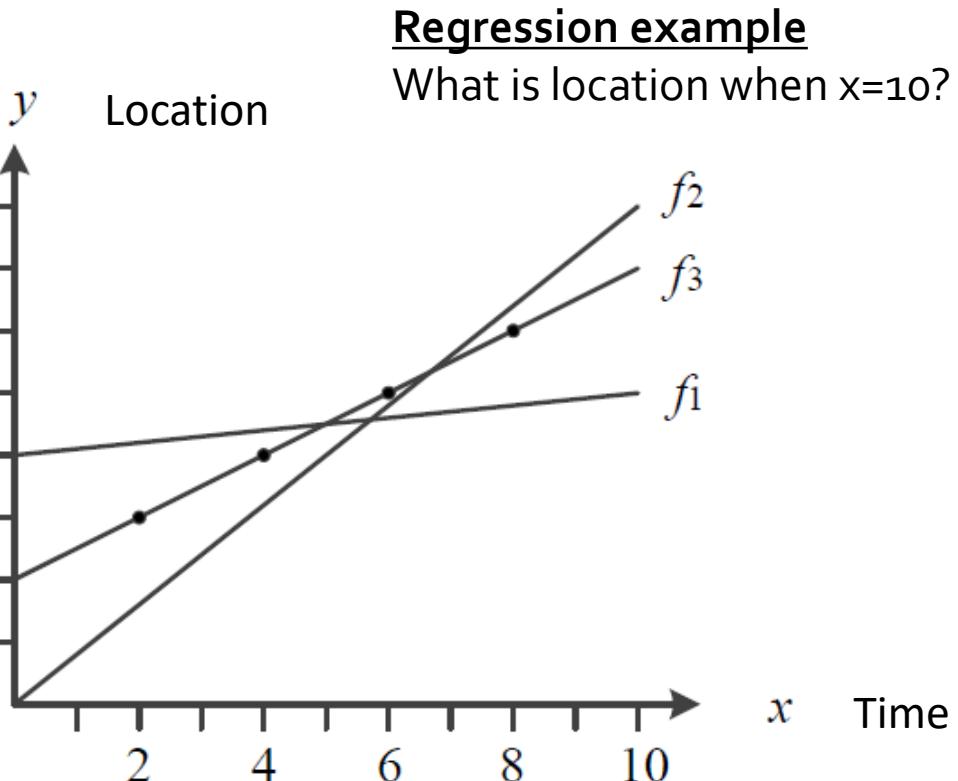
Play video games

$$f(\text{[video game screenshot]}) = \text{“move left”}$$

Machine Learning: Basic Concept

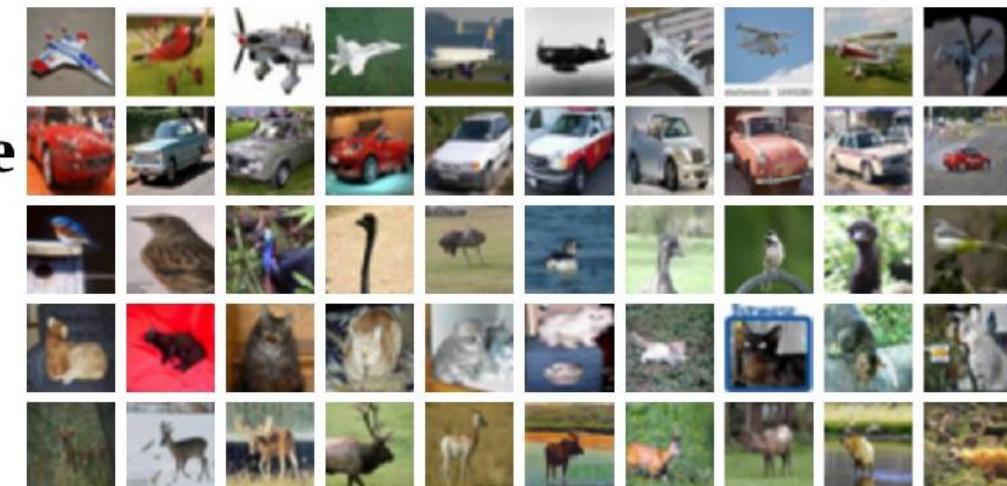
- **Prediction task**

- Regression: returns a specific value
- Classification: returns a class label



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

airplane



automobile



bird



cat



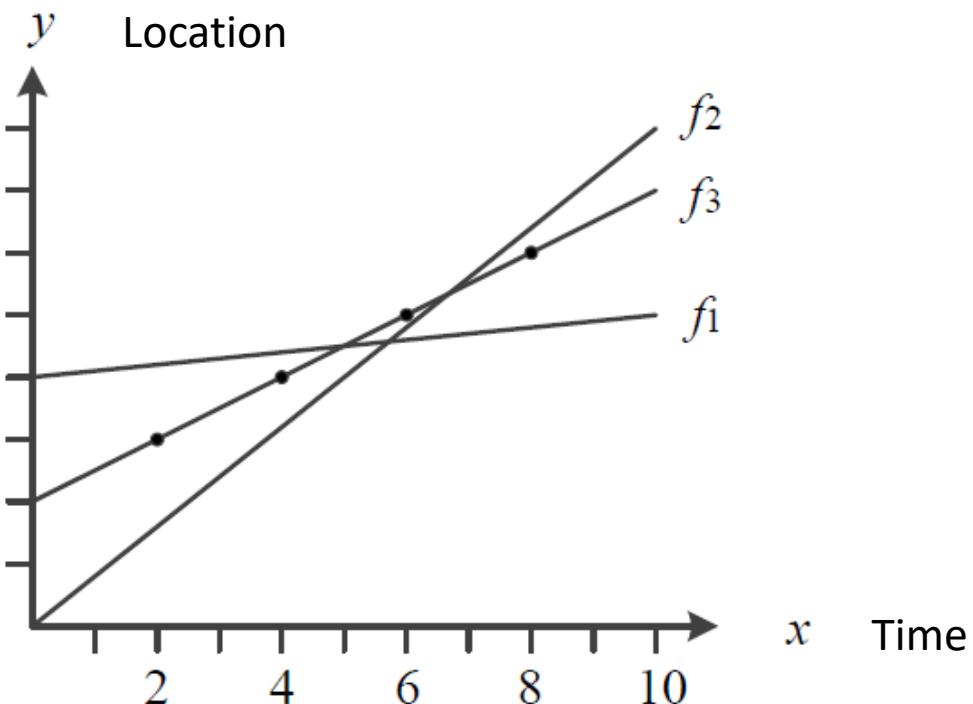
deer



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

$$y = wx + b$$

Optimal parameter $w=0.5$ $b=2.0$

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: $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$ 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

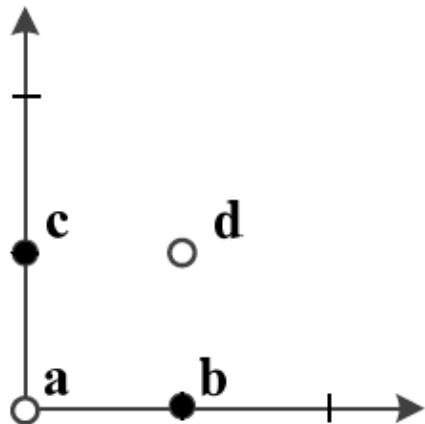
$$y = \underline{w_1}x_1 + \underline{w_2}x_2 + \dots + \underline{w_d}x_d + \underline{b}$$

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

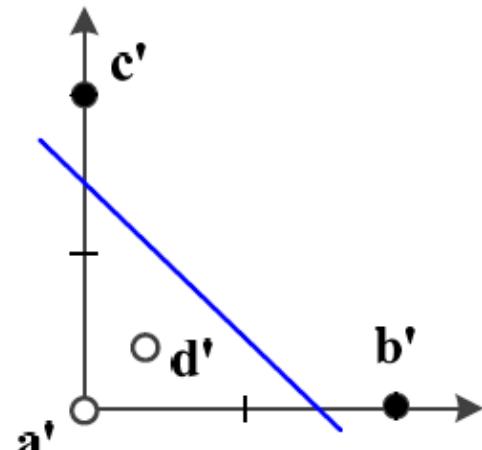
$$y = \underline{w_1}x_1^2 + \underline{w_2}x_2^2 + \dots + \underline{w_d}x_d^2 + \underline{w_{\frac{d(d+1)}{2}}}x_1x_2 + \dots + \underline{w_{\frac{d(d+1)}{2}}}x_{d-1}x_d + \underline{w_{\frac{d(d+1)}{2}+1}}x_1 \dots + \underline{w_{\frac{d(d+1)}{2}+d}}x_d + \underline{b}$$

Machine Learning: Basic Concept

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



Original feature space



Transformed feature space,
which is linearly separable

Toy example

$$\mathbf{x} = (x_1, x_2)^T \rightarrow \mathbf{x}' = \left(\frac{x_1}{2x_1x_2 + 0.5}, \frac{x_2}{2x_1x_2 + 0.5} \right)^T$$

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

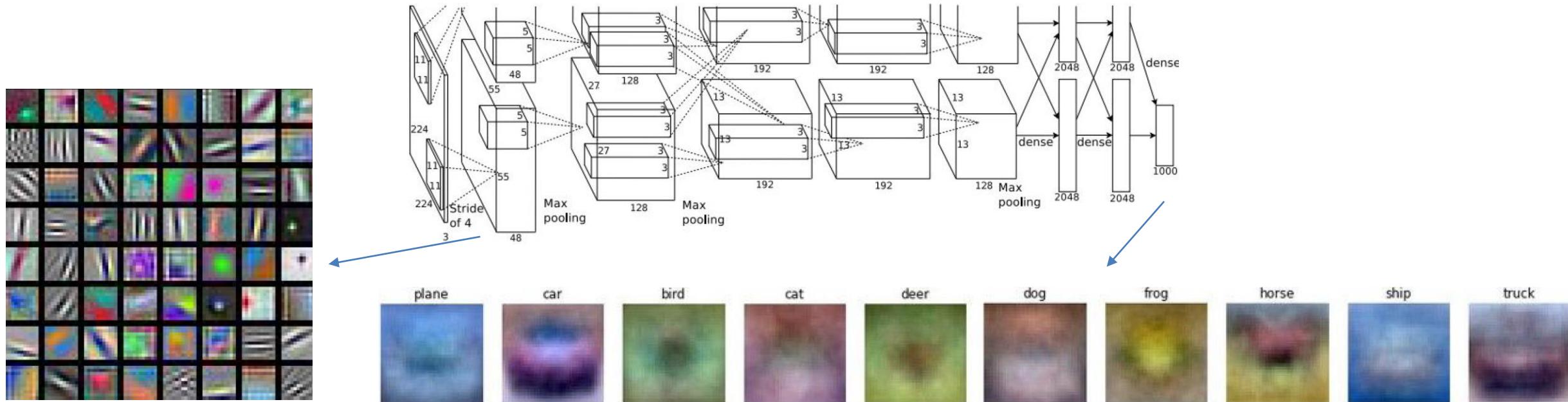
$$\mathbf{b} = (1,0)^T \rightarrow \mathbf{b}' = (2,0)^T$$

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

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

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.)

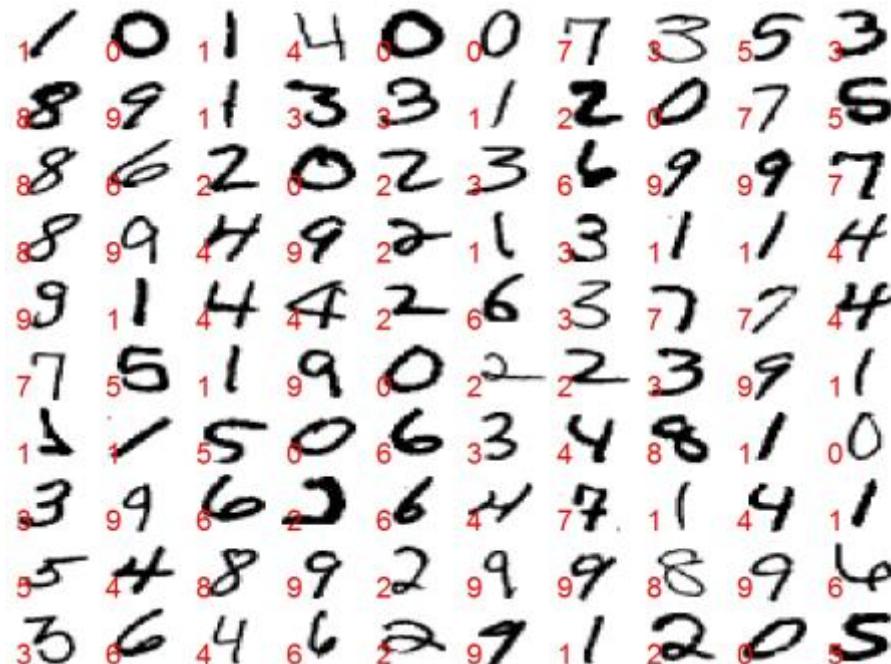


Data for Machine Learning

- **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.

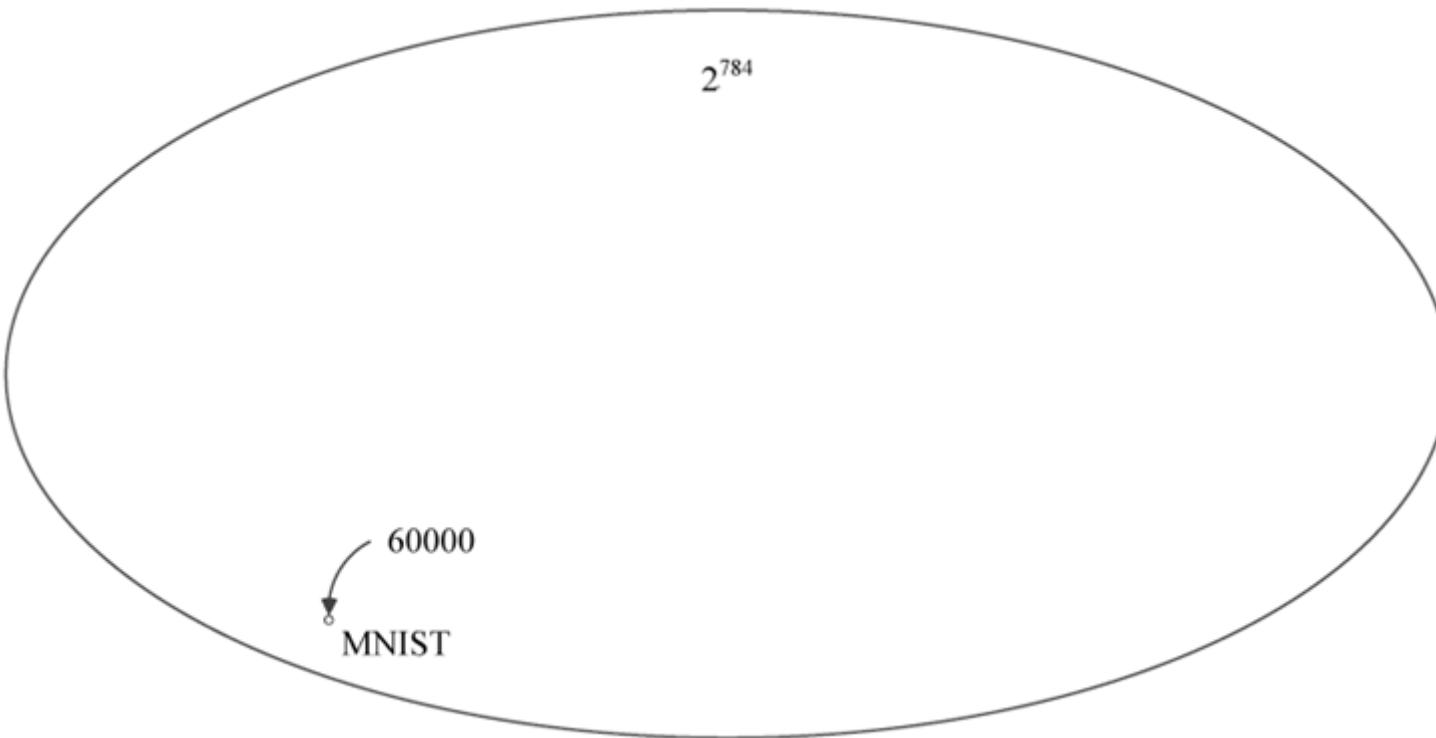
- **MNIST database**

- 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.

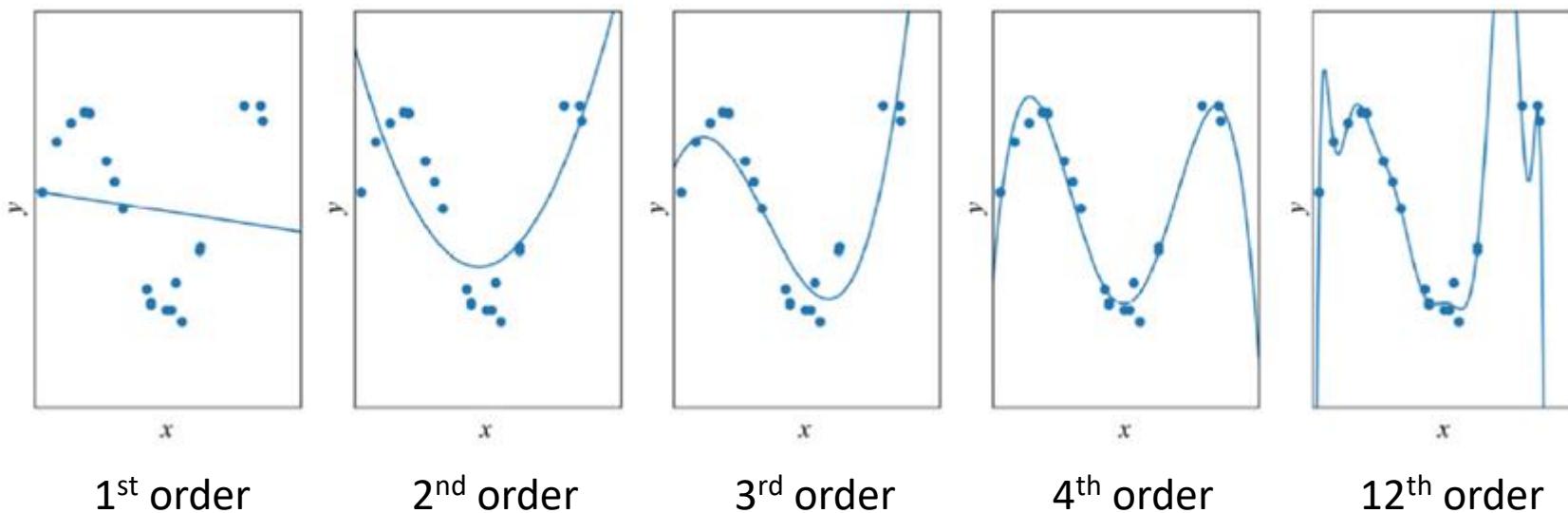


Data for Machine Learning

- **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.
 - Manifold assumption
Smooth change according to certain rules like 

Training Model: under-fitting vs. over-fitting

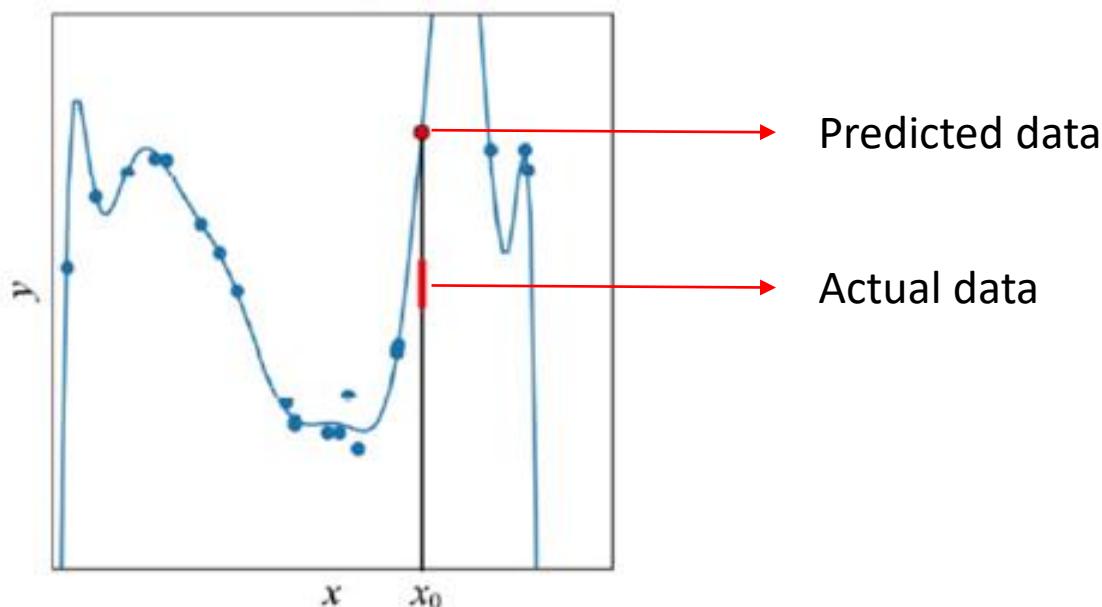
- **Under-fitting**
 - Model capacity is too small to fit the data accordingly.
 - Model with higher order can be used.



Training Model: under-fitting vs. over-fitting

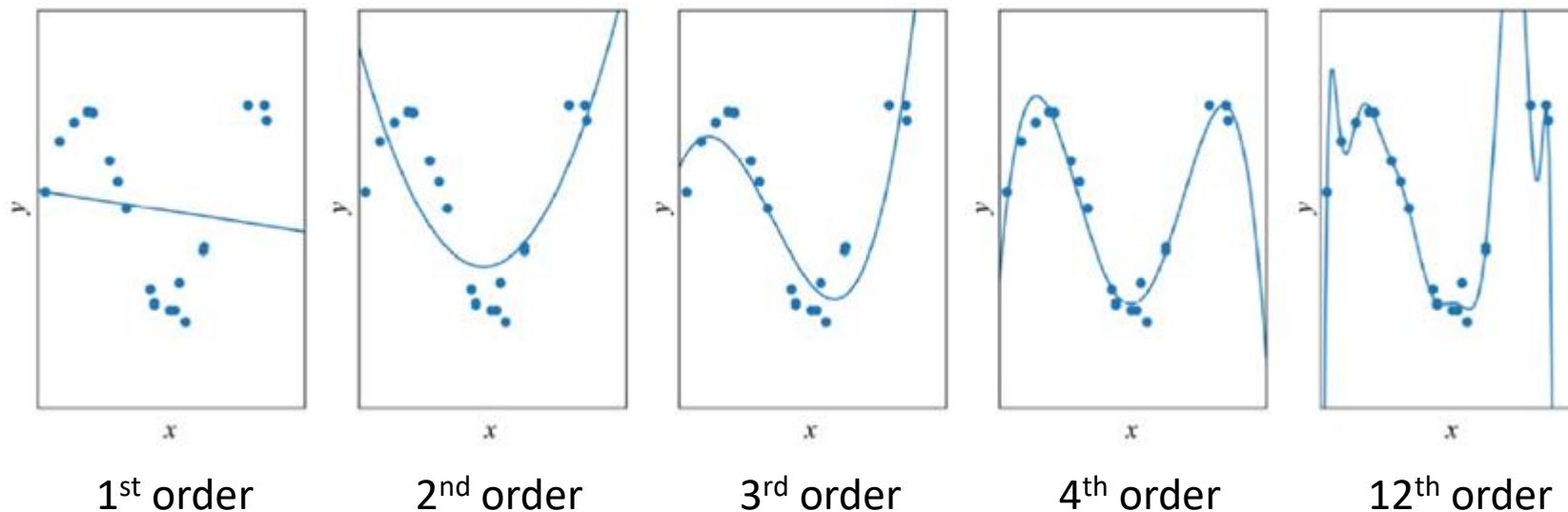
- **Over-fitting**

- 12th 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.

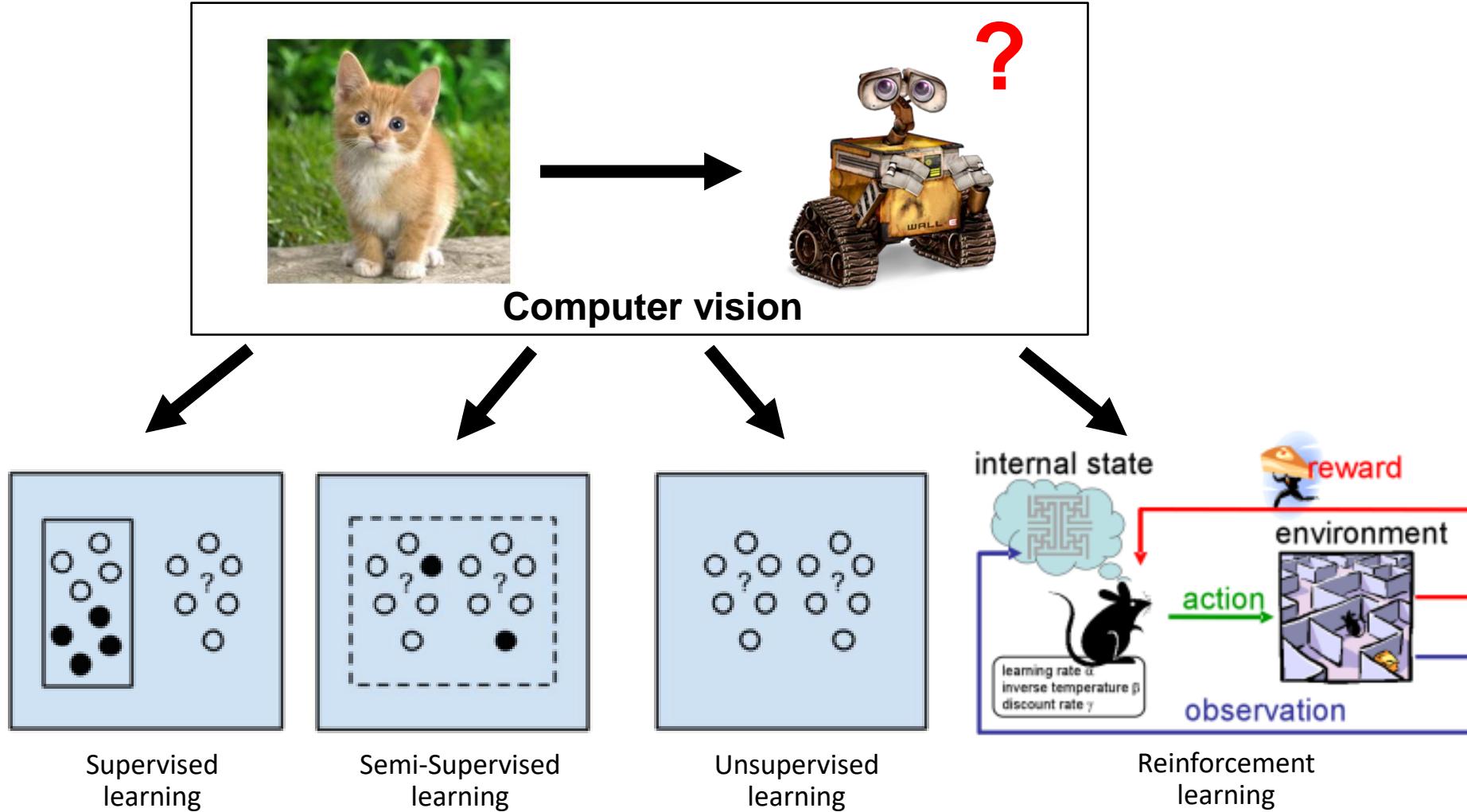


Training Model: under-fitting vs. over-fitting

- 1st and 2nd order model show poor performance for both the training and the test set.
- 12th order model shows high performance in training set, but low performance in test set. → low generalization ability
- 3rd and 4th order model are lower than the 12th order model for the training set, but the test set has high performance → higher generalization capability



Spectrum of supervision



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.

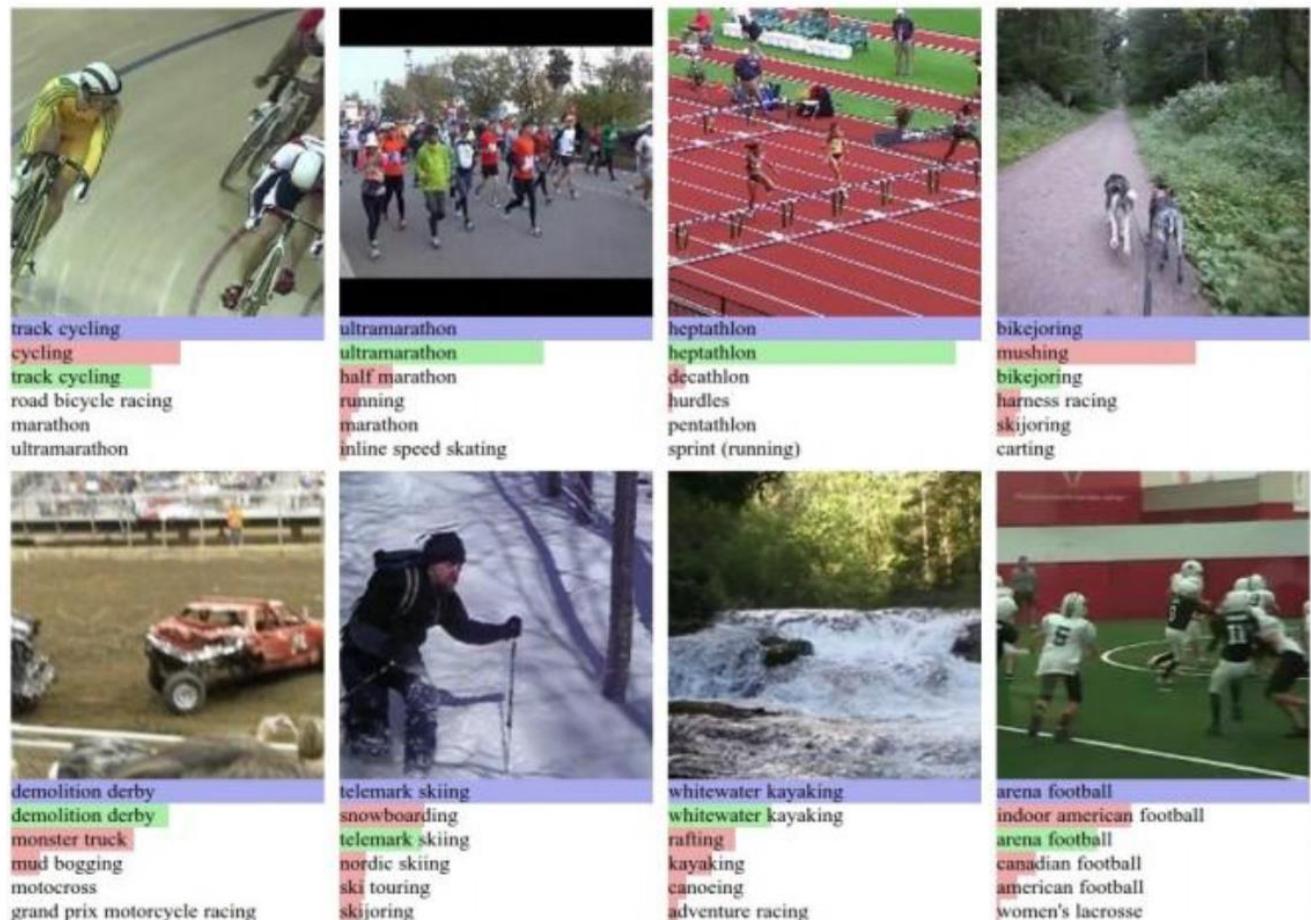
Introduction to Deep Learning

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.

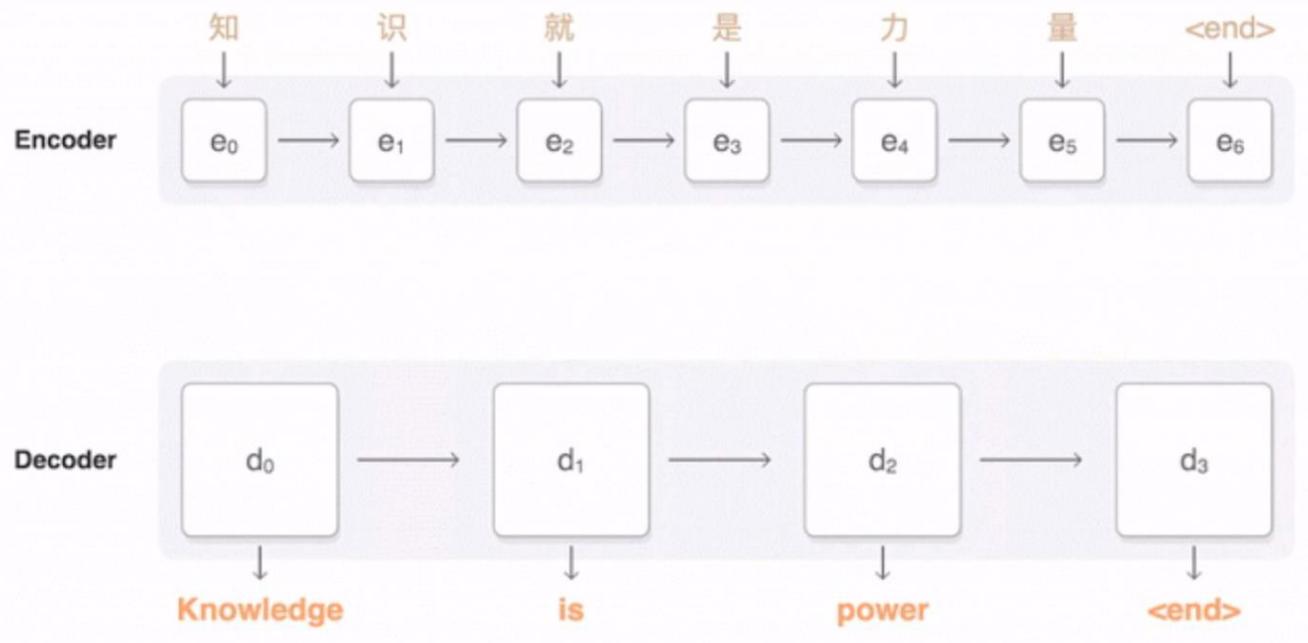
Deep learning success

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



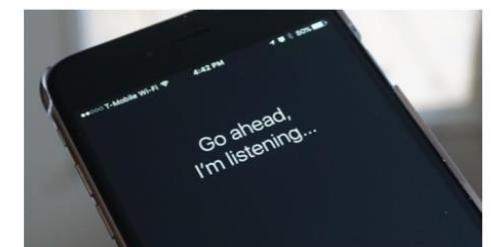
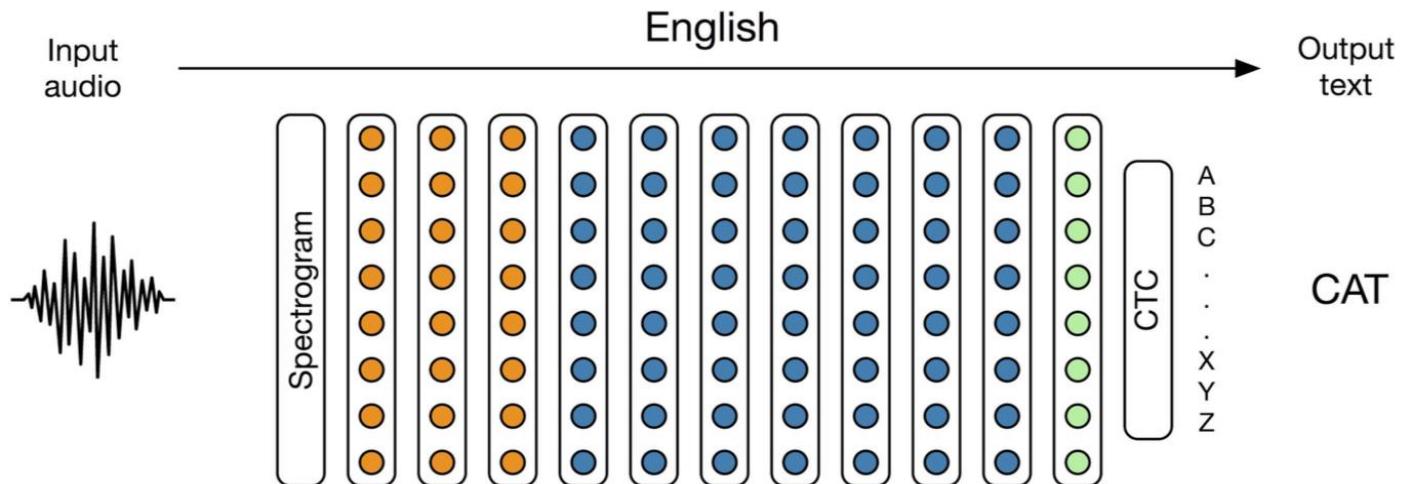
Deep learning success

- 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



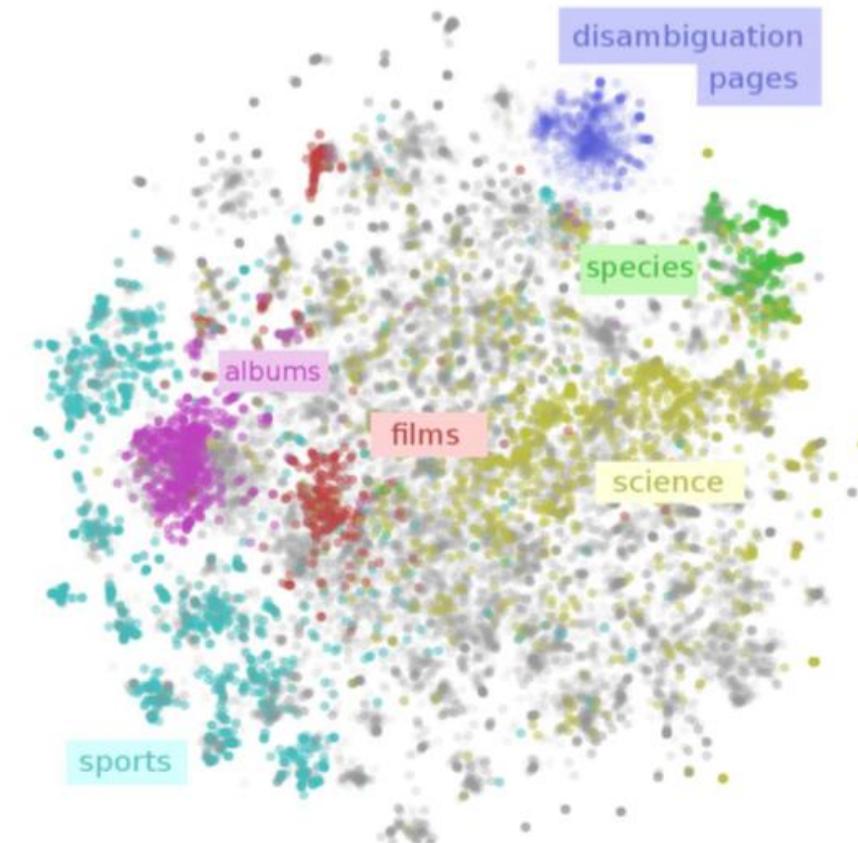
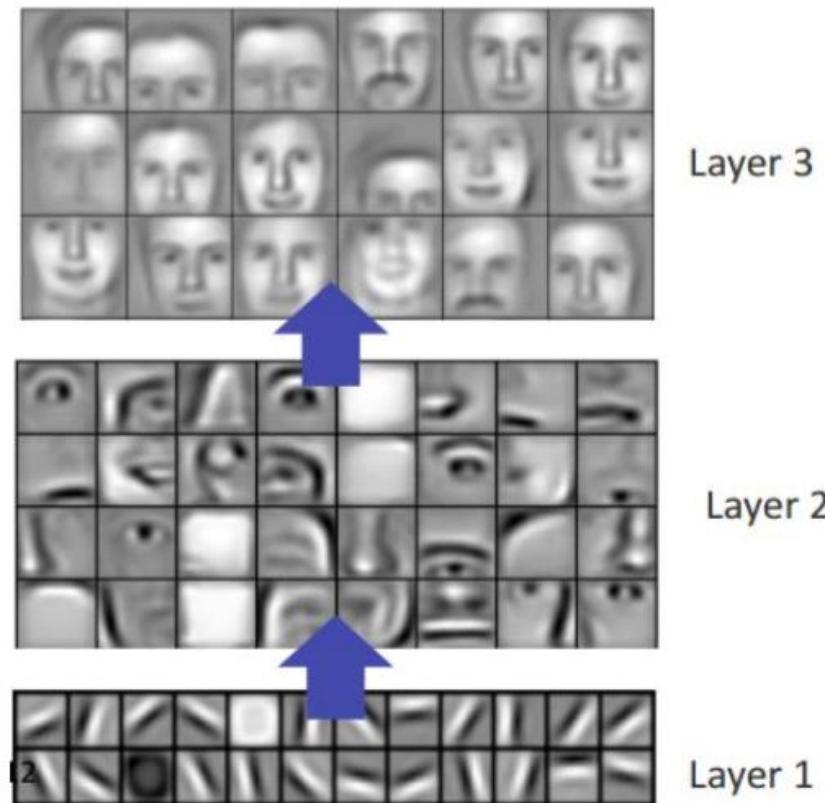
Deep learning success

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



Why deep learning?

- Hand-crafted features vs. Learned features



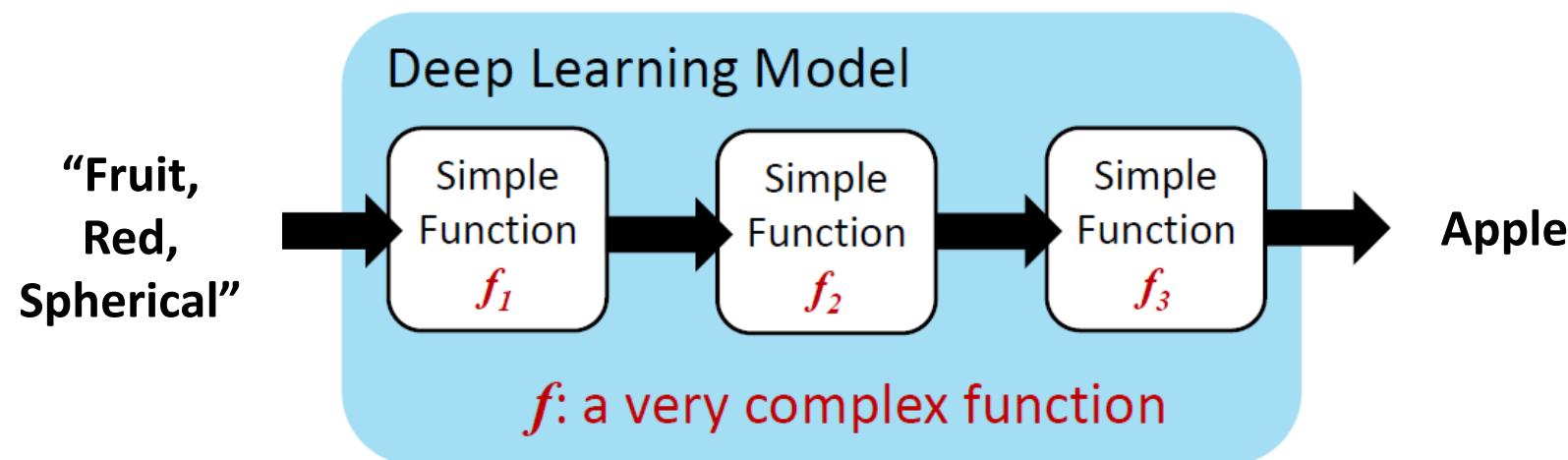
Why now?

- Large datasets
- GPU hardware advances + Price decreases
- Improved techniques (algorithm)



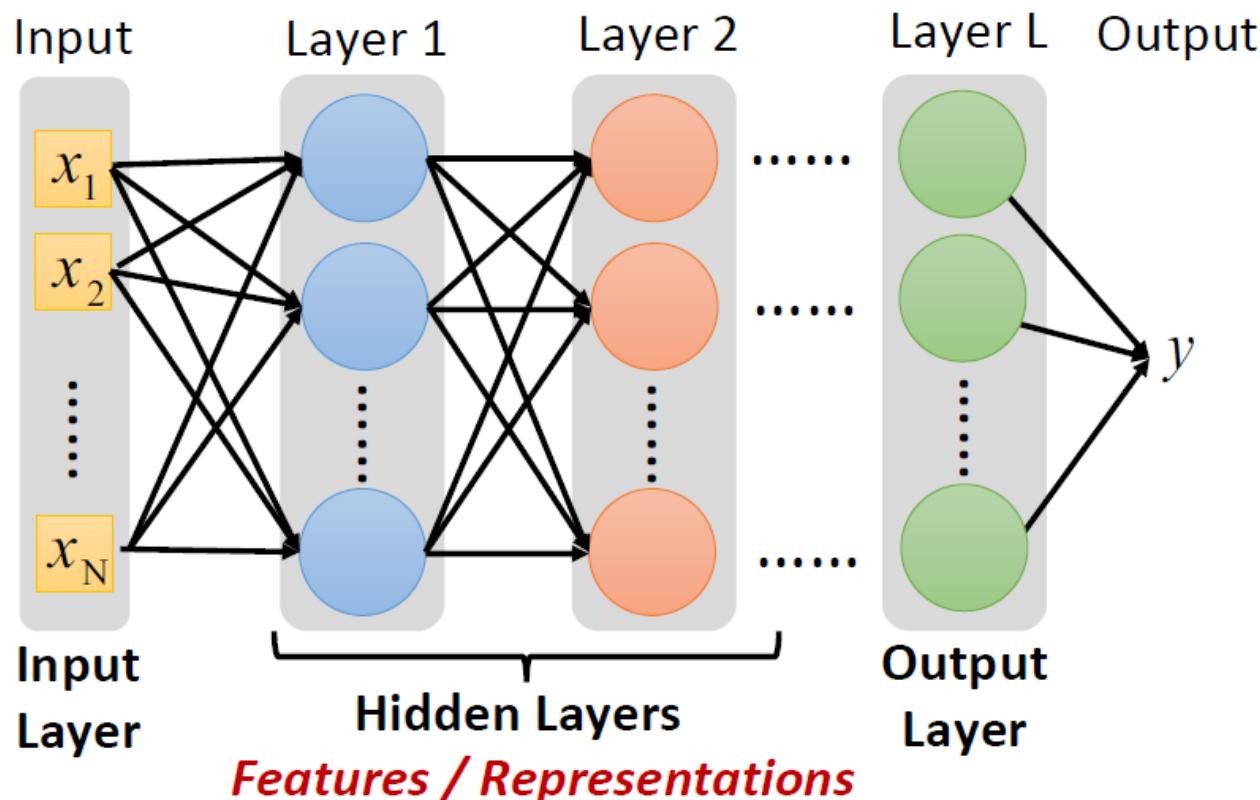
What is Deep Learning?

- **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



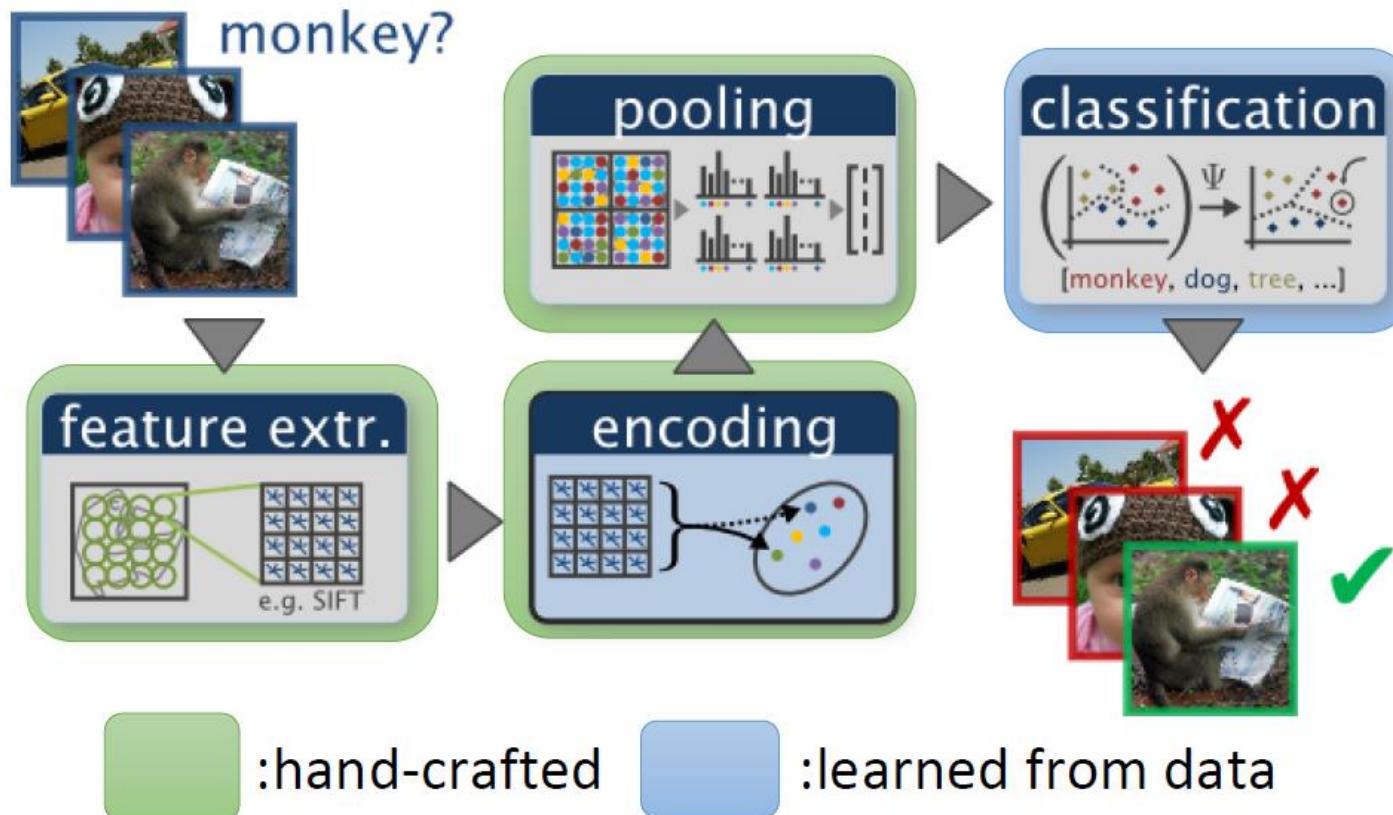
What is Deep Learning?

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



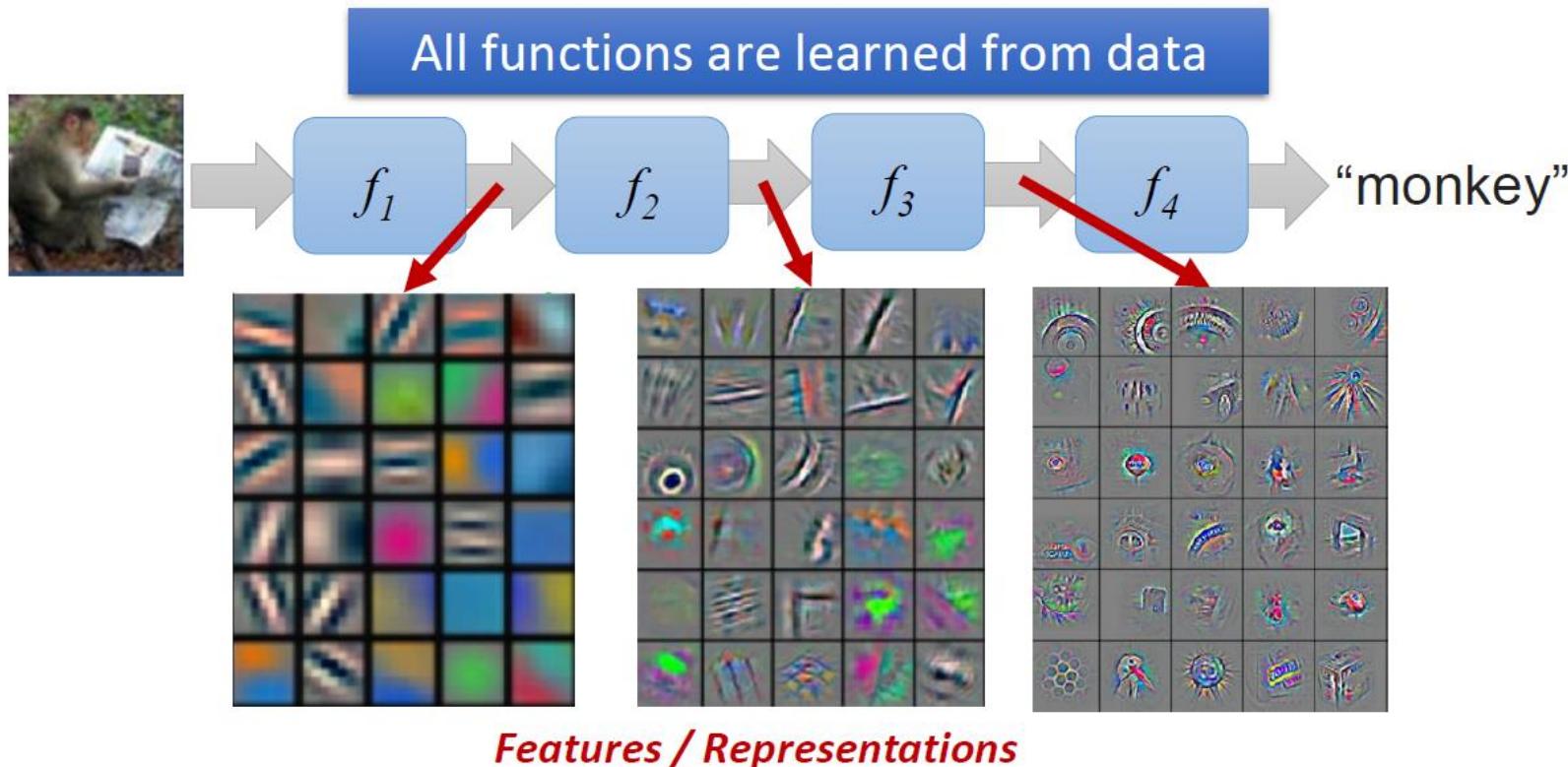
Machine Learning vs. Deep Learning

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



Machine Learning vs. Deep Learning

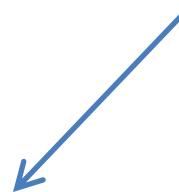
- Deep vs Shallow: Image Recognition
 - Deep model using deep learning



Machine Learning vs. Deep Learning

- Machine Learning vs. Deep Learning

Machine Learning = Feature descriptor + Classifier



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

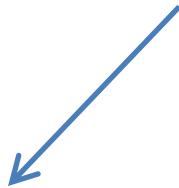


Optimizing the classifier weights on features
Ex) Nearest Neighbor (NN), Support Vector Machine
(SVM), Random Forest (RF)

Machine Learning vs. Deep Learning

- Machine Learning vs. Deep Learning

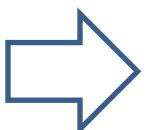
Deep Learning = Feature descriptor + Classifier



Feature: Representation learned by machine
Automatically learned internal knowledge



Optimizing the classifier weights on
features

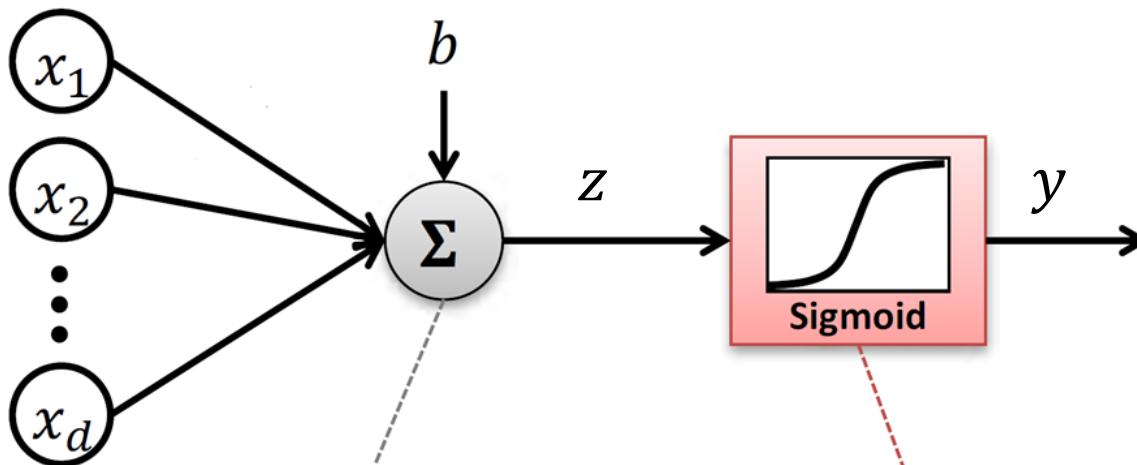


Neural network based model

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

Deep Learning

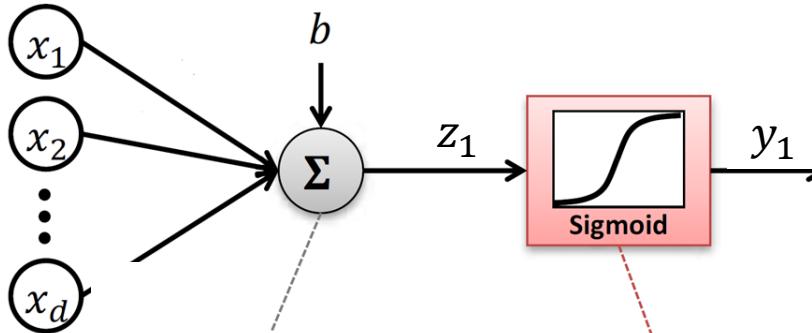
- A single neuron



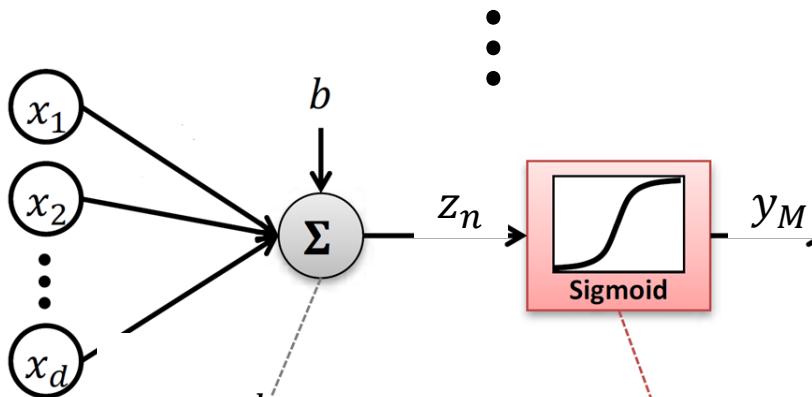
$$z = \mathbf{w}^T \mathbf{x} + b = \sum_{k=1}^d w_k x_k + b \quad y = \frac{1}{1 + e^{-z}}$$

Deep Learning

- A single layer with multiple neurons



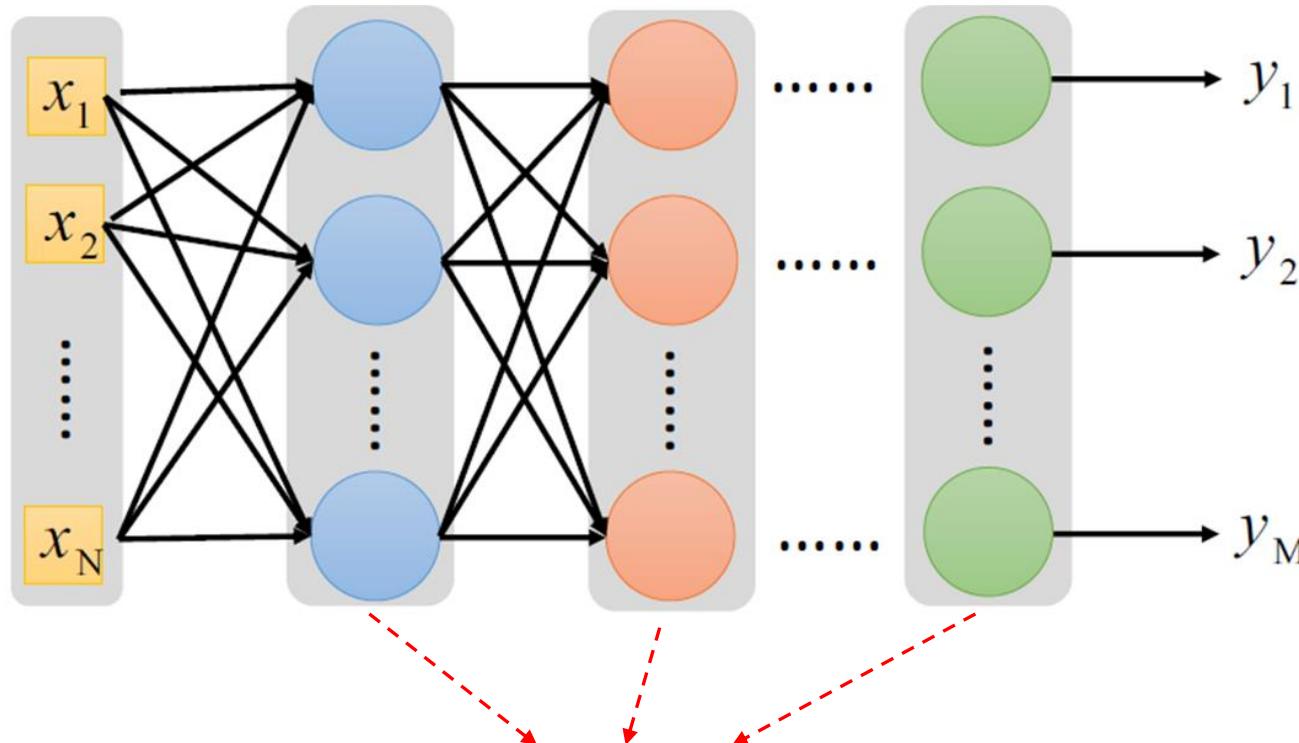
$$z_1 = \mathbf{w}_1^T \mathbf{x} + b_1 = \sum_{k=1}^d w_{1k} x_k + b_1 \quad y_1 = \frac{1}{1 + e^{-z_1}}$$



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

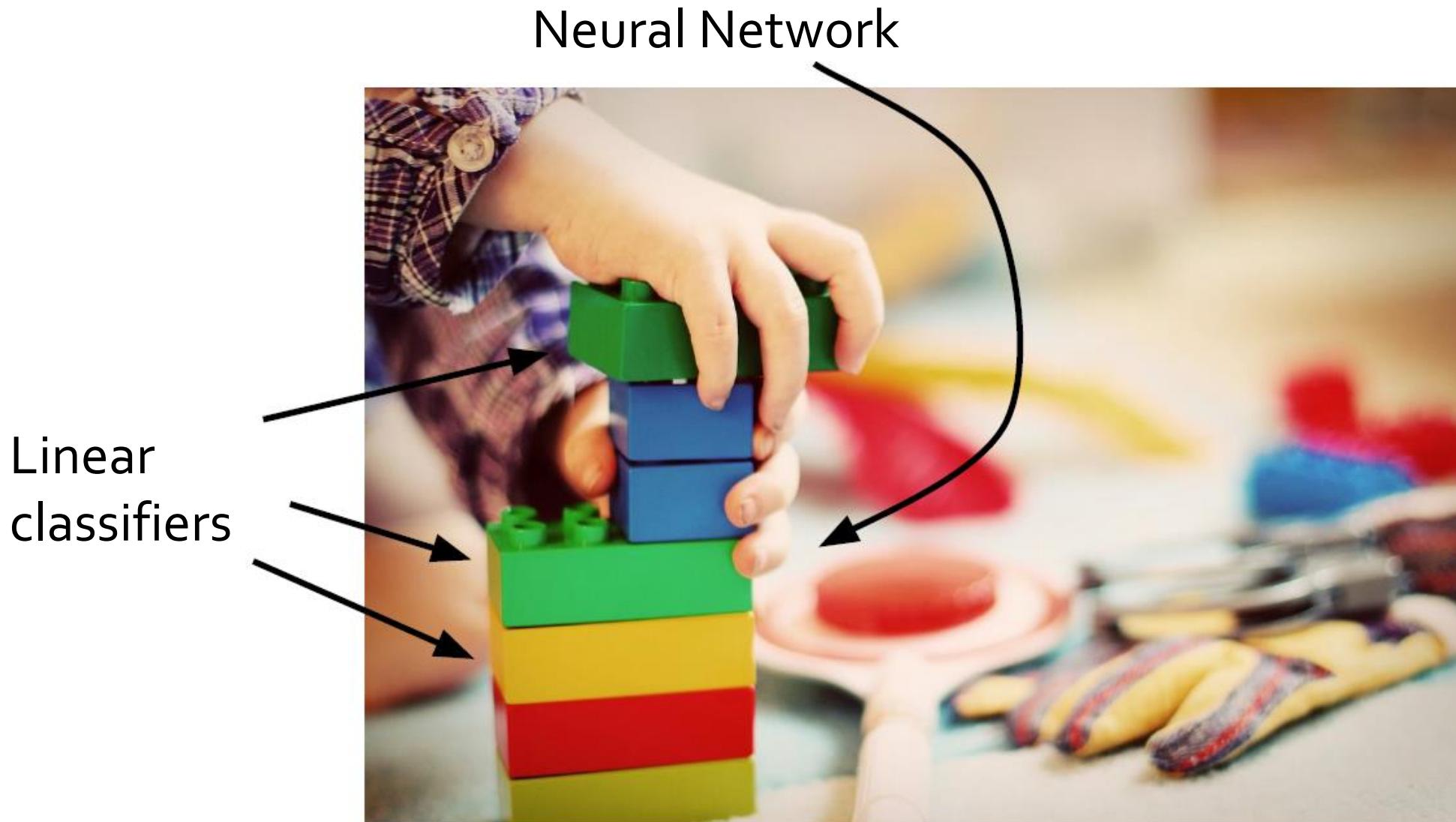
Deep Learning

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



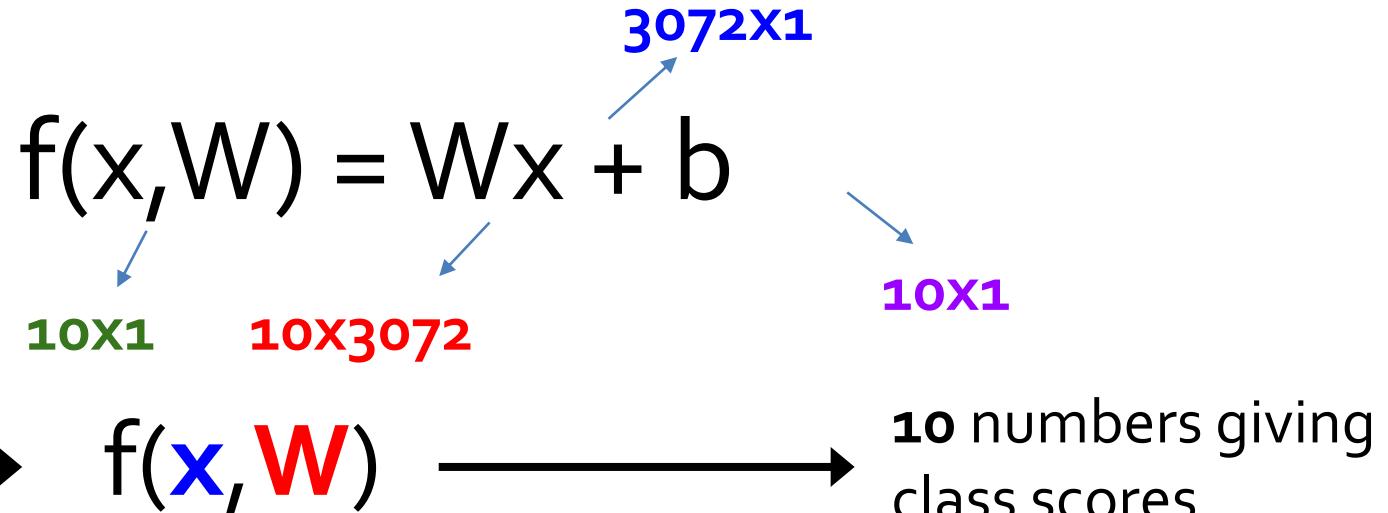
Each layer consists of the linear classifier
and activation function

Linear classifier



Parametric Approach

- (Review) Unit 3 ML basics and classification



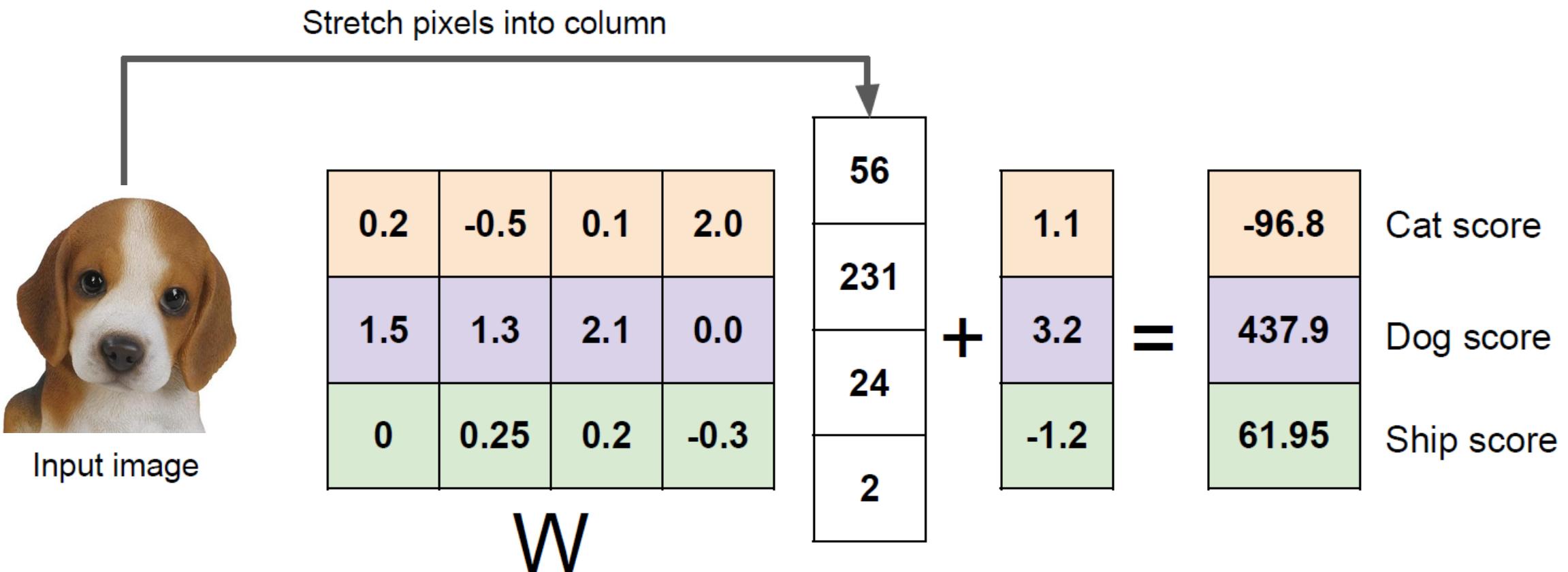
Array of $32 \times 32 \times 3$ numbers
(3072 numbers total)

W
parameters
(or weights)

Parametric Approach

- (Review) Unit 3 recognition

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



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

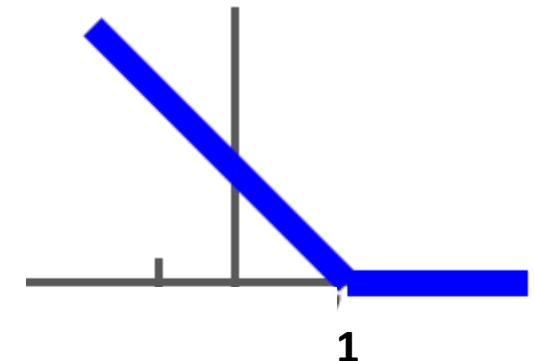
- **Loss function**
 - quantifies our unhappiness with the scores across the training data.
- **Type of loss function**
 - Hinge loss
 - Cross-entropy loss
 - Log likelihood loss
 - Regression loss

Loss Function: Hinge Loss

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

$$L_i = \max(0, 1 - y_i \cdot s) \quad s = \mathbf{w}^T \mathbf{x}_i + b$$

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



- **Hinge loss (=multiclass SVM loss)**

- C : The number of class (> 2)

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

\mathbf{x}_i : input data (e.g. image)

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

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

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

Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.

With some \mathbf{W} the scores $f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ are



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

Given a dataset of examples

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

\mathbf{x}_i : image

y_i : class label (integer)

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

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

Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.

With some \mathbf{W} the scores $f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ are



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

Multiclass SVM loss (=hinge loss)

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

where score vector $\mathbf{s} = f(\mathbf{x}_i, \mathbf{W})$

Loss Function: Hinge Loss

Suppose: 3 training examples, 3 classes.

With some \mathbf{W} the scores $f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ are



| | | | |
|------|------------|------------|-------------|
| 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 |

Multiclass SVM loss (=hinge loss)

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

$$\begin{aligned} &= \max(0, 5.1 - 3.2 + 1) + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 = 2.9 \end{aligned}$$

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 = 0 \end{aligned}$$

$$\begin{aligned} &= \max(0, 2.2 - (-3.1) + 1) + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 6.3) + \max(0, 6.6) \\ &= 6.3 + 6.6 \end{aligned}$$

Loss over full dataset is average

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

Loss Function: Log Likelihood Loss

- **Log likelihood loss**

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

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

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

\mathbf{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 $C = 10$.

$$\mathbf{z}_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \quad \Rightarrow \quad L_i = -\log 0.7$$

Loss Function: Cross-entropy Loss

- Cross-entropy loss

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

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

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

\mathbf{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 $C = 10$.

$$\mathbf{z}_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$
$$\mathbf{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix}$$


$$L_i = -\log(1 - 0.1) - \log 0.7 - \log(1 - 0.2)$$

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



unnormalized probabilities

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

Softmax activation
function

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

x_i : image

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

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

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

unnormalized log probabilities

probabilities

Softmax + Log Likelihood Loss



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

unnormalized probabilities

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

$$L_i = -\log(0.13) \\ = 0.89$$

unnormalized log probabilities

probabilities

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

Softmax + Cross-entropy Loss



unnormalized probabilities

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

cat
car
frog

| |
|------|
| 3.2 |
| 5.1 |
| -1.7 |

exp

| |
|-------|
| 24.5 |
| 164.0 |
| 0.18 |

normalize

| |
|------|
| 0.13 |
| 0.87 |
| 0.00 |

unnormalized log probabilities

probabilities

$$\begin{aligned} L_i &= -\log(0.13) - \log(1-0.87) - \log(1-0.0) \\ &= 0.89 * 2 \end{aligned}$$

Loss Function: Regression Loss

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

$$L_i = |y_i - s_i|$$

$$L_i = (y_i - s_i)^2$$

$$y_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad s_i = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix} \quad \Rightarrow \quad L_i = |y_i - s_i| = |0 - 0.1| + |1 - 0.7| + |0 - 0.2|$$

Regularization

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

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:

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 = 0 \end{aligned}$$

With W twice as large:

$$\begin{aligned} &= \max(0, 2.6 - 9.8 + 1) + \max(0, 4.0 - 9.8 + 1) \\ &= \max(0, -6.2) + \max(0, -4.8) \\ &= 0 + 0 = 0 \end{aligned}$$

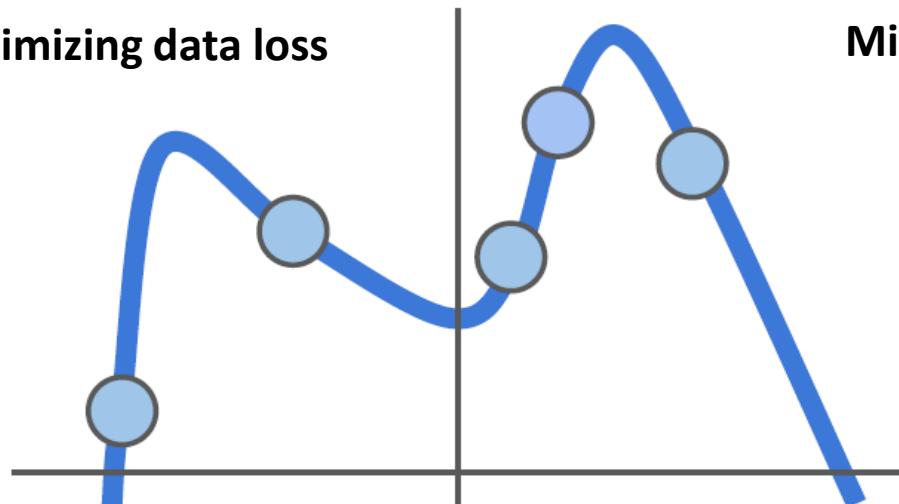
Regularization

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

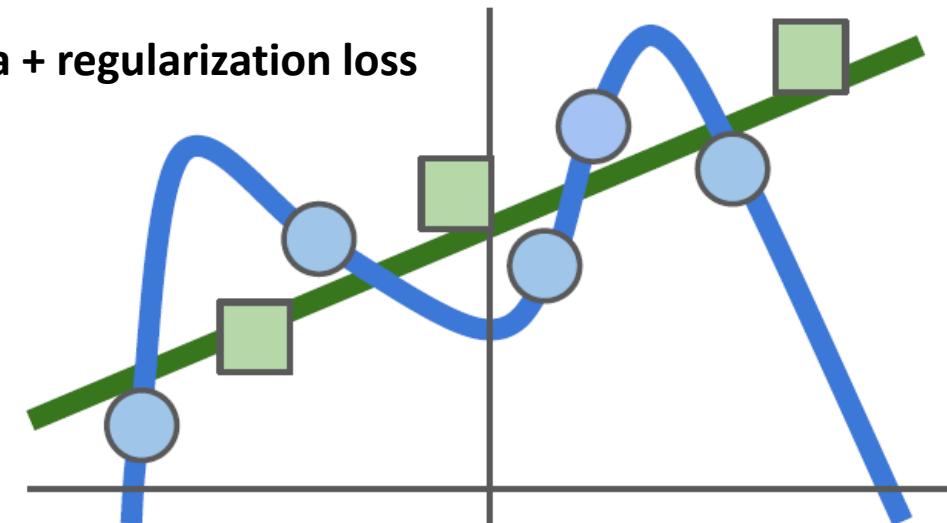
Data loss: Model predictions should match training data

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

Minimizing data loss



Minimizing data + regularization loss



Regularization

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

λ : regularization strength
(hyperparameter)

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

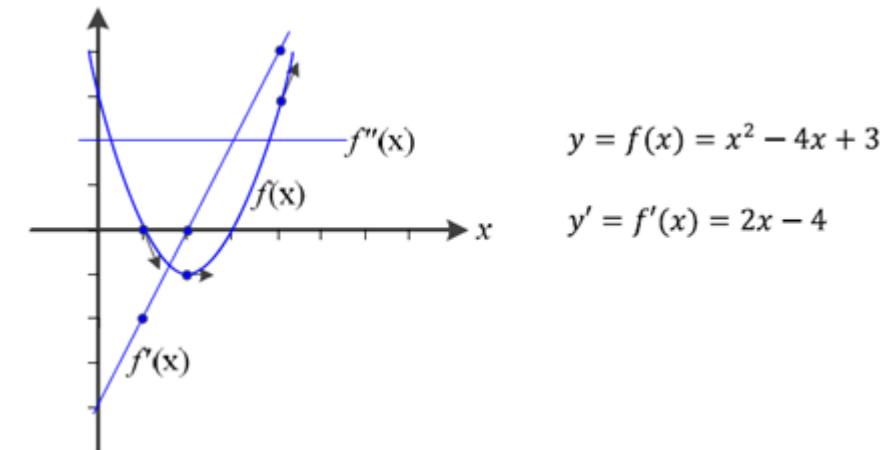
- L₂ regularization: $R(\mathbf{W}) = \sum_{k,l} W_{k,l}^2$
- L₁ regularization: $R(\mathbf{W}) = \sum_{k,l} |W_{k,l}|$
- Elastic net (L₁ + L₂): $R(\mathbf{W}) = \sum_{k,l} \beta W_{k,l}^2 + |W_{k,l}|$
- Max norm regularization: $|\mathbf{w}_j^T| < c$ **for all j**
- Dropout (will see later)
- Batch normalization, stochastic depth (will see later)

Optimization: Gradient Descent

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

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

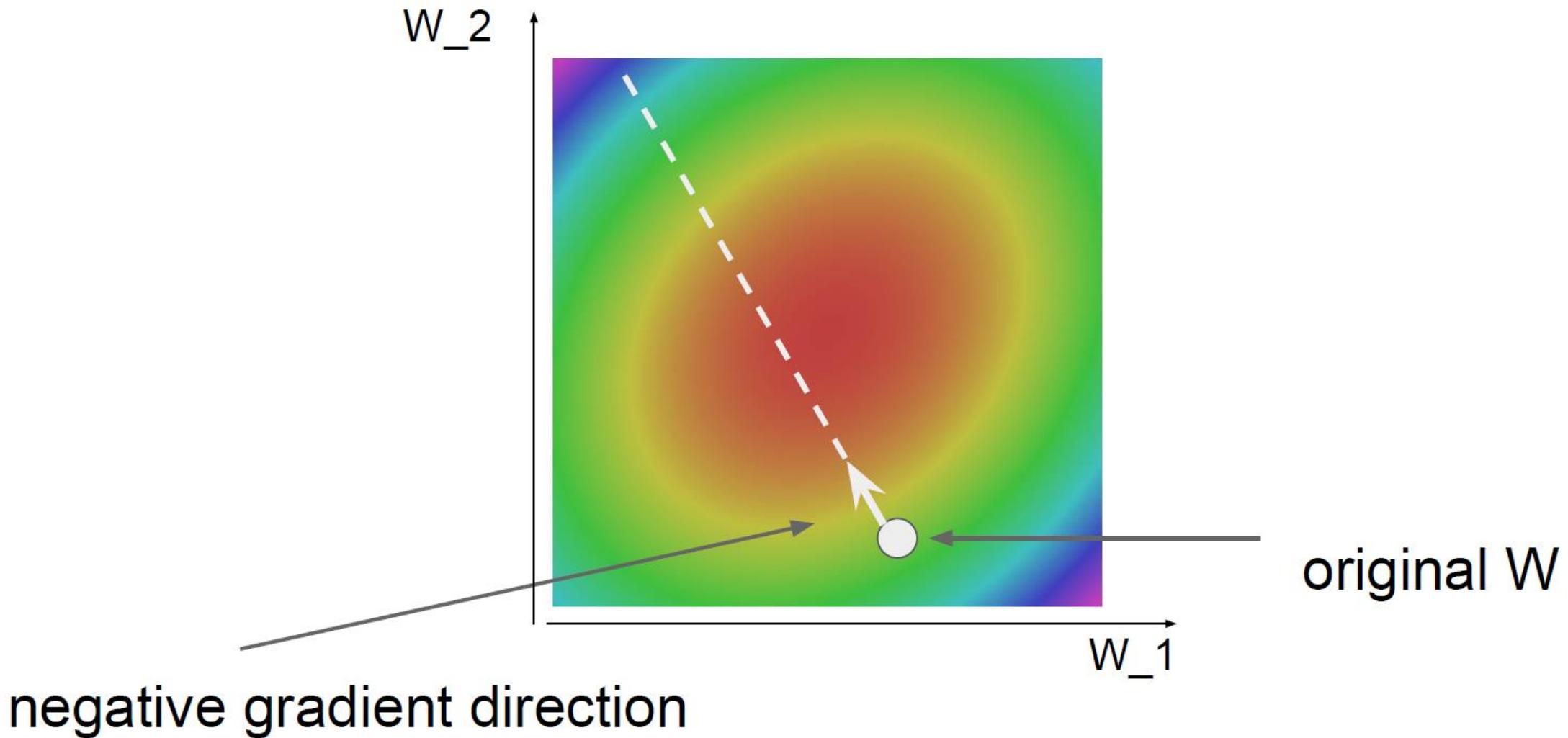
- α : step size (a.k.a. learning rate)



```
# Vanilla Gradient Descent

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

Optimization: Gradient Descent



Optimization: Stochastic Gradient Descent (SGD)

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

Full sum is too expensive when N is large!

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

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

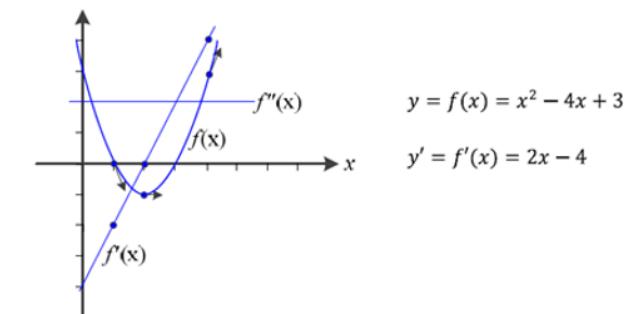
Backpropagation

- A widely used algorithm for training feedforward neural networks.
- A way of computing gradients of expressions through recursive application of chain rule.
 - Backpropagation computes the gradient of the loss function with respect to the weights of the network (model) for a single input–output example.

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

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

- α : step size (a.k.a. learning rate)



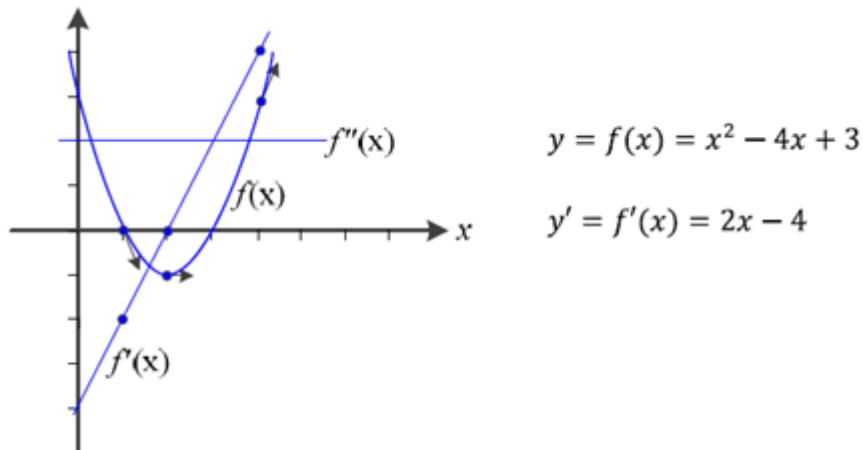
Derivative

- Optimization using derivative

- 1st order derivative

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

- $f'(x)$: 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: $-f'(x)$



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

Derivative

- **Partial derivative**

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

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

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

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

Derivative

- **Jacobian matrix**

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

$$\mathbf{J} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_d} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_d} \end{pmatrix}$$

Ex) $\mathbf{f}: \mathbb{R}^2 \mapsto \mathbb{R}^3 \quad \mathbf{f}(\mathbf{x}) = (2x_1 + x_2^2, -x_1^2 + 3x_2, 4x_1x_2)^T$

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

- **Hessian matrix**

- 2nd order partial derivative matrix

$$\mathbf{H} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 x_1} & \frac{\partial^2 f}{\partial x_1 x_2} & \dots & \frac{\partial^2 f}{\partial x_1 x_n} \\ \frac{\partial^2 f}{\partial x_2 x_1} & \frac{\partial^2 f}{\partial x_2 x_2} & \dots & \frac{\partial^2 f}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n x_1} & \frac{\partial^2 f}{\partial x_n x_2} & \dots & \frac{\partial^2 f}{\partial x_n x_n} \end{pmatrix}$$

Ex) $f(\mathbf{x}) = f(x_1, x_2)$
 $= \left(4 - 2.1x_1^2 + \frac{x_1^4}{3}\right)x_1^2 + x_1x_2 + (-4 + 4x_2^2)x_2^2$

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

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

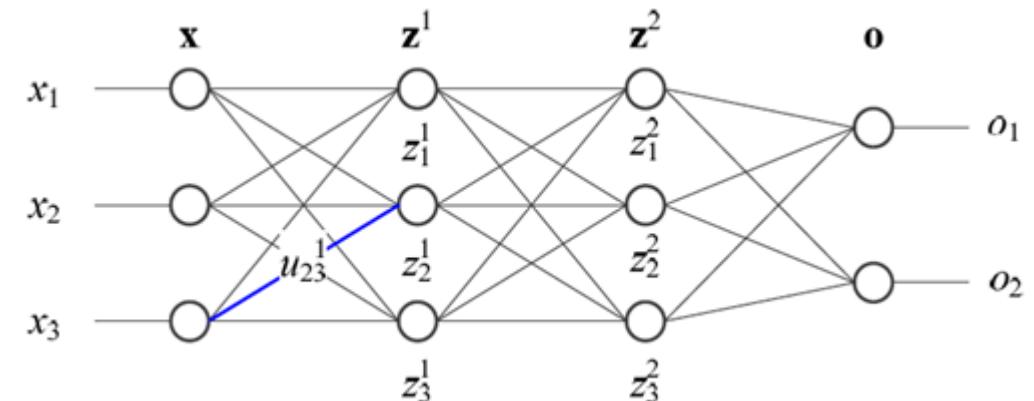
Derivative

- Chain rule

$$\begin{array}{ccc} f(x) = g(h(x)) & \xrightarrow{\hspace{1cm}} & f'(x) = g'(h(x))h'(x) \\ f(x) = g(h(i(x))) & & f'(x) = g'(h(i(x)))h'(i(x))i'(x) \end{array}$$

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

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

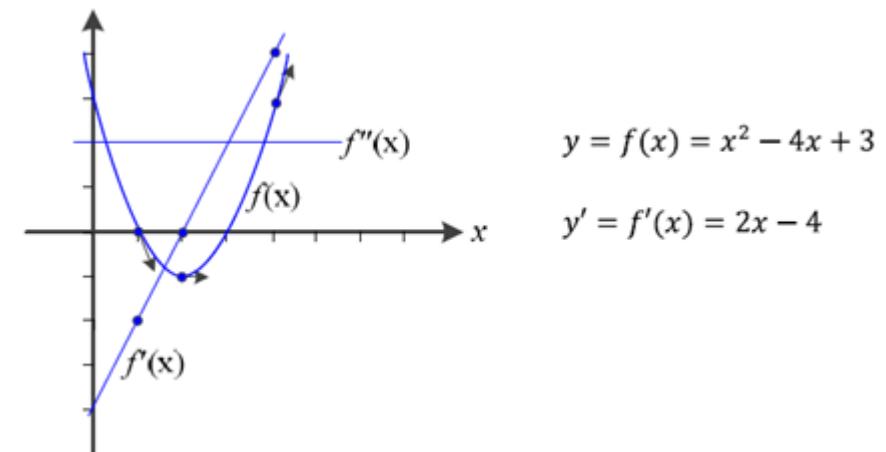


Why are we talking about derivatives?

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

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

- α : step size (a.k.a. learning rate)



```
# Vanilla Gradient Descent

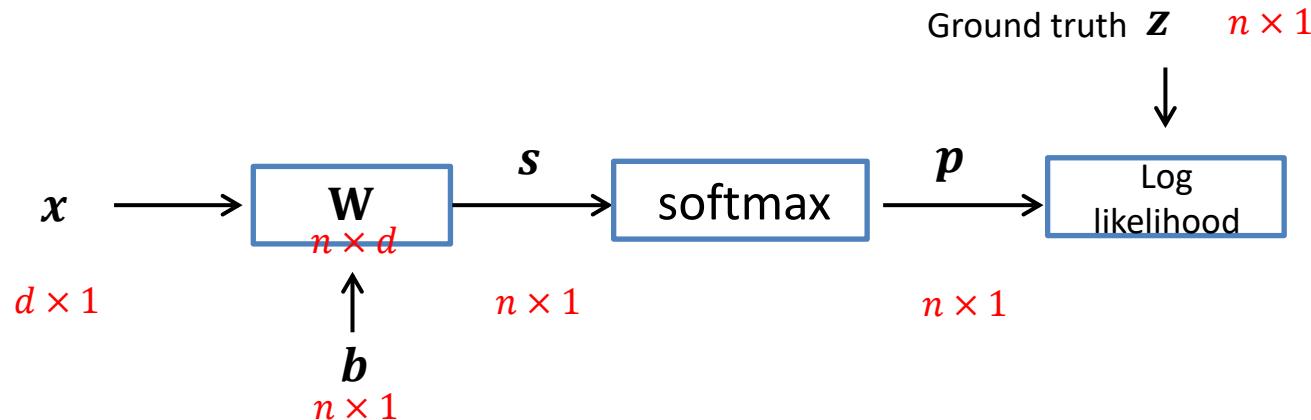
while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

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 $\frac{\partial L}{\partial \mathbf{b}}$



$$\frac{\partial L}{\partial \mathbf{p}} \quad \frac{\partial L}{\partial s} = \frac{\partial p}{\partial s} \frac{\partial L}{\partial p} \quad \frac{\partial L}{\partial \mathbf{w}_j} = \frac{\partial s}{\partial \mathbf{w}_j} \frac{\partial L}{\partial s} \quad \frac{\partial L}{\partial \mathbf{W}} = \left(\frac{\partial L}{\partial \mathbf{w}_1} \quad \frac{\partial L}{\partial \mathbf{w}_2} \quad \dots \quad \frac{\partial L}{\partial \mathbf{w}_n} \right)^T$$
$$\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial s}{\partial \mathbf{b}} \frac{\partial L}{\partial s} = \frac{\partial L}{\partial s}$$

Analytic Gradient: Linear Equation

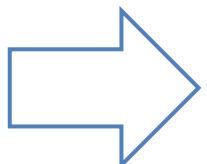
$$s = \mathbf{W}\mathbf{x} + \mathbf{b} \quad \longleftrightarrow \quad \begin{aligned} s_1 &= \mathbf{w}_1^T \mathbf{x} + b_1 \\ s_2 &= \mathbf{w}_2^T \mathbf{x} + b_2 \\ &\vdots \\ s_n &= \mathbf{w}_n^T \mathbf{x} + b_n \end{aligned}$$
$$\mathbf{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}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

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

$$\frac{\partial s_2}{\partial \mathbf{w}_1} = \mathbf{0}$$

\vdots

$$\frac{\partial s_n}{\partial \mathbf{w}_1} = \mathbf{0}$$



$$\frac{\partial \mathbf{s}}{\partial \mathbf{w}_1} = [\mathbf{x} \ 0 \ 0 \ \cdots \ 0] \in \Re^{d \times n}$$

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

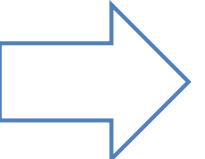
jth column

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

Analytic Gradient: Linear Equation

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

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

$$\frac{\partial s_2}{\partial \mathbf{x}} = \mathbf{w}_2$$

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

\vdots

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

Analytic Gradient: Sigmoid Function

- **Sigmoid function**

For a scalar x

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

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

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

for $j = 1, \dots, n$

Analytic Gradient: Softmax Activation Function

- Softmax function

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

score function $\mathbf{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}$ probability $\mathbf{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{pmatrix}$

- 1st order derivative of softmax function

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

Analytic Gradient: Softmax Activation Function

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

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

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

- **1st order derivative of binary hinge loss**

$$L = \max(0, 1 - y \cdot s)$$

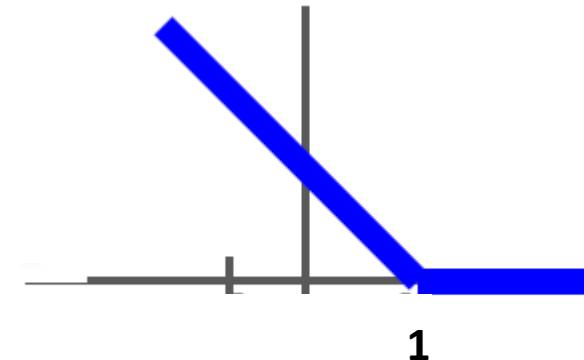
$$= \begin{cases} 1 - y \cdot s & \text{if } 1 - y \cdot s > 0 \\ 0 & \text{otherwise} \end{cases}$$



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

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

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



Analytic Gradient: Hinge Loss

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

- **1st order derivative of hinge loss**

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

x : image

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

$$\mathbf{s} = \mathbf{Wx} + \mathbf{b}$$

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



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

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

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

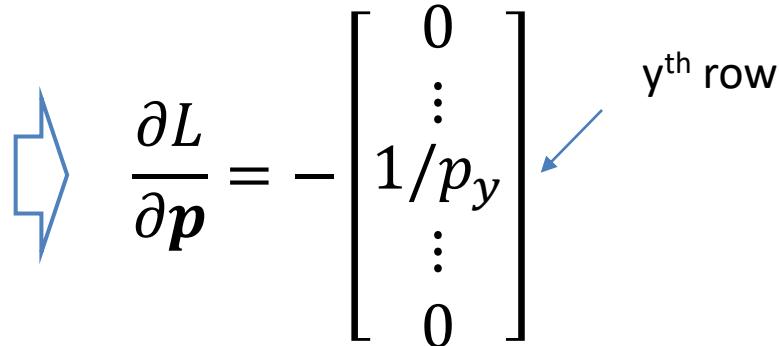
Analytic Gradient: Log Likelihood Loss

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

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

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

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


$$\frac{\partial L}{\partial \mathbf{p}} = - \begin{bmatrix} 0 \\ \vdots \\ 1/p_y \\ \vdots \\ 0 \end{bmatrix}$$

y^{th} row

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

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

$\mathbf{z} = (z_1 z_2 \dots z_n)^T$, $z_y = 1$ and $z_{k \neq y} = 0$

Example

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

$$\mathbf{z} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\mathbf{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix}$$



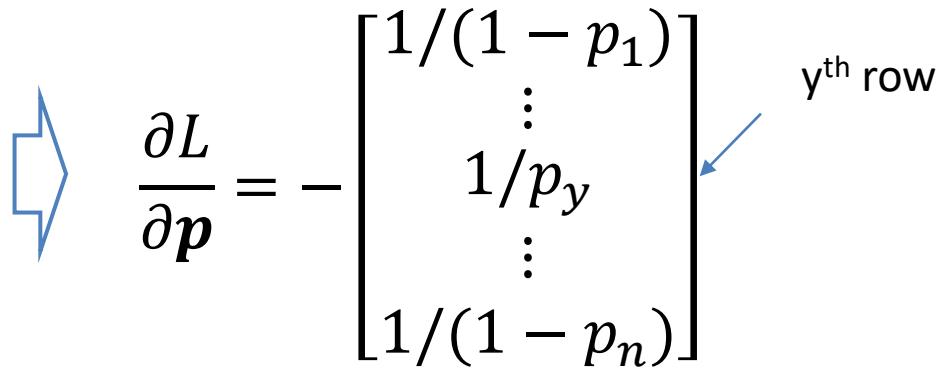
$$\frac{\partial L}{\partial \mathbf{p}} = - \begin{bmatrix} 0 \\ 1/0.7 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Analytic Gradient: Cross-entropy Loss

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

$$L = - \sum_{j=1}^n (z_j \log p_j + (1 - z_j) \log(1 - p_j)) \quad \mathbf{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{pmatrix}$$

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


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

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

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

$\mathbf{z} = (z_1 \ z_2 \ \dots \ z_n)^T$, $z_y = 1$ and $z_{k \neq y} = 0$

Example

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

$$\mathbf{z} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{p} = \begin{pmatrix} 0.1 \\ 0.7 \\ 0 \\ \vdots \\ 0.2 \end{pmatrix}$$

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

Analytic Gradient: Regression Loss

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

- Regression loss

$$L = (\mathbf{y} - \mathbf{s})^2$$



$$\frac{\partial L}{\partial \mathbf{s}} = -2(\mathbf{y} - \mathbf{s})$$

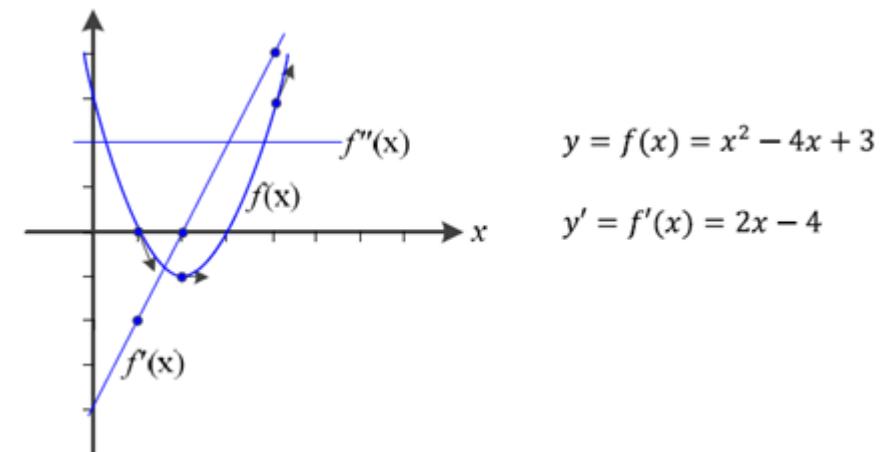
$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{s} = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix}$$

Why are we talking about derivatives?

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

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

- α : 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
```