# Grasping Deep Learning from Fundamentals to Applications

*June 15, 2023*
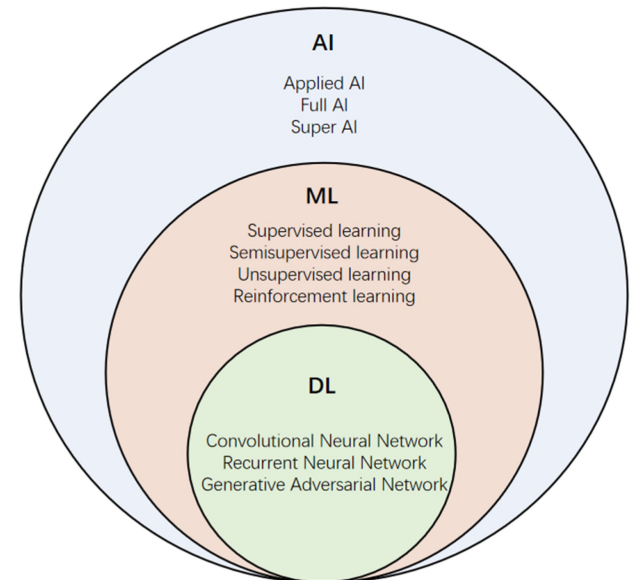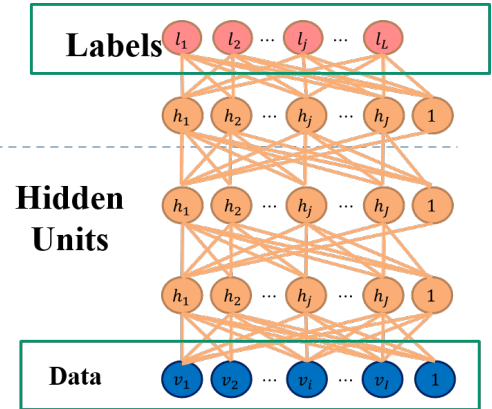
## Lecture 1 – Intro to Deep Learning

Instructors: **Yufei Huang**, PhD; **Arun Das**, PhD
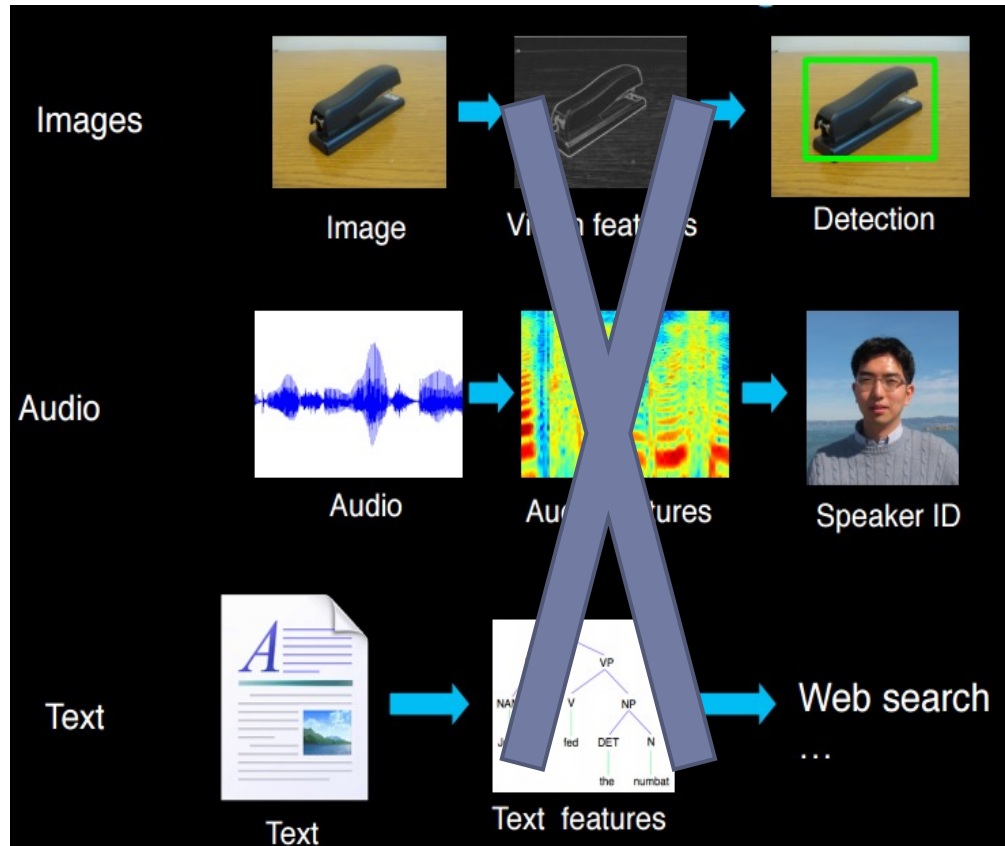
University of Pittsburgh

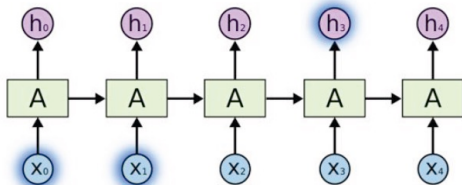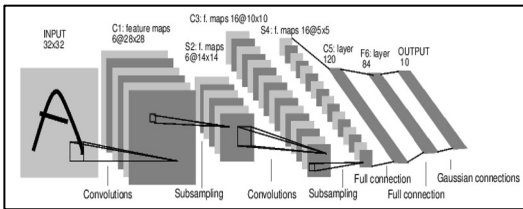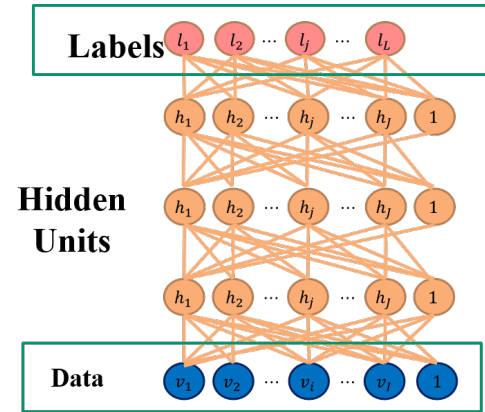# What is deep learning?

What is deep learning?
- **A model paradigm**
  - Neural networks
- **A learning paradigm**
  - Supervised
  - Unsupervised
  - Reinforcement
  - Self-supervised
  - Transfer
  - Contrastive
  - …





University of Pittsburgh

# Deep Learning brings deep revolution for AI

## Simplify AI systems
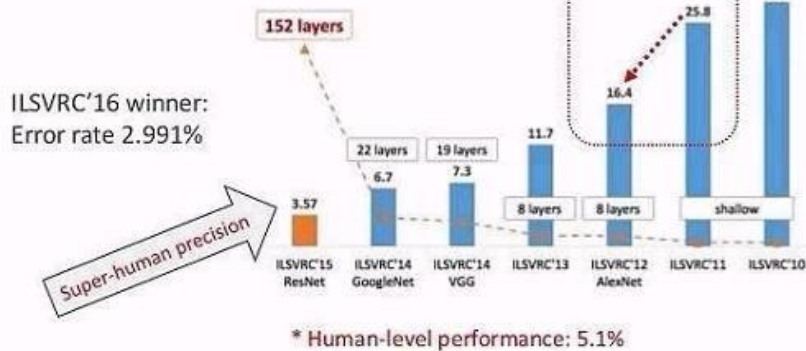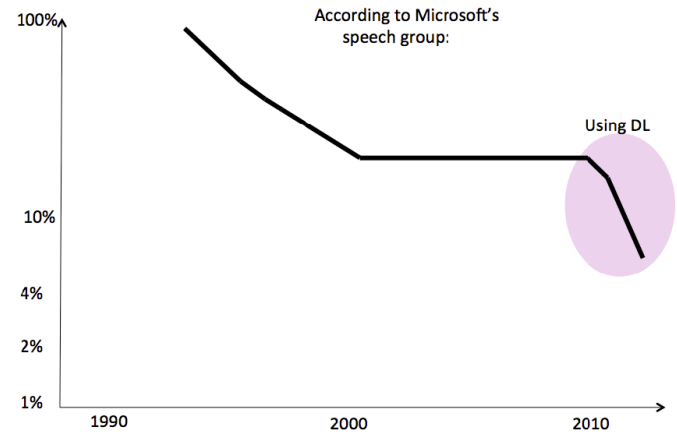
# Deep Learning brings deep revolution for AI

## Maximize AI performance

Computer vision



Speech recognition



Deep Learning in Speech Recognition

**Microsoft achieves 'human parity' in speech recognition system**

University of Pittsburgh

# History of DL



**1958** Perceptron

**1974** Backpropagation

Convolution Neural Networks for Handwritten Recognition
**1998**

Google Brain Project on 16k Cores
**2012**

awkward silence (AI Winter)

**1969** Perceptron criticized

**1995** SVM reigns

**2006** Restricted Boltzmann Machine

**2012** AlexNet wins ImageNet
IMAGENET

University of Pittsburgh

# History of DL – Artificial Neuron, 1943



## McCulloch and Pitts
"A Logical Calculus of the Ideas Immanent in Nervous Activity"

# History of DL – Perceptron, 1958



1957  Frank Rosenblatt



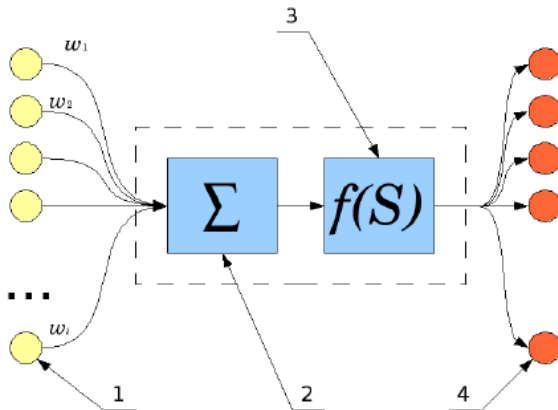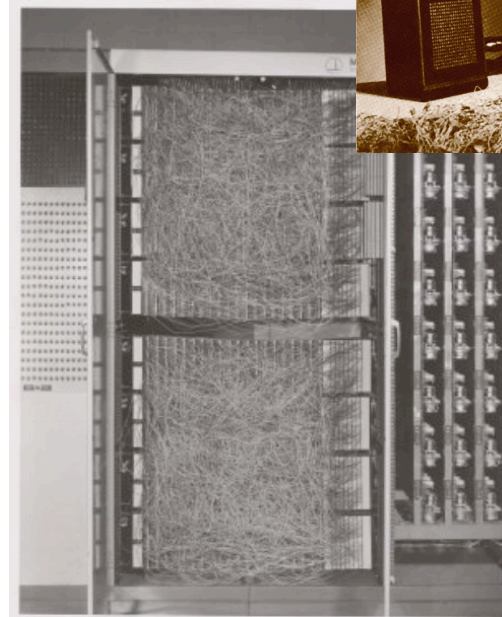*"[The Perceptron is] the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."*

THE NEW YORK TIMES

University of Pittsburgh

# History of DL – Backpropagation, 1974



➡ Measure how **small changes** in weights **affect** output

➡ Can apply NN to **regression**

(1974)     1986

**(Werbos)**  Rumelhart, Hinton, Williams

"Learning representations by back-propagating errors" (Nature)

➡ **Multilayer** neural networks, etc.

# History of DL –
# NN winter; late 1990 – early 2000

- Not enough data

- Not enough computing power

- Imperfect activation function

University of
Pittsburgh

# History of DL –
# Restricted Boltzmann machine 2005



Geoffrey Hinton "The Godfather of deep learning"



98.2% on the MNIST test set

# History of DL – Convolution 1998, 2010

Yann LeCun, NYU & Meta



99.50% on the MNIST test set
CURRENT BEST: 99.77% by committee of 35 conv. nets

# History of DL – ImageNet 2012 -



**Convolution layers**    **Fully connected layers**

**Image recognition**
- **AlexNet (2012)**
  - 7 layers;
  - 1000 labels; >1.2M images
  - 17% vs 25.7%
- **GoogLeNet; VGG (2014)**
  - 19 layers (VGG)
  - ~7%

- **Deep Residual Net (2015)**
  - 152 layers
  - 3.57% > human recognition

University of Pittsburgh

ImageNet Challenge

ILSVRC'16 winner:
Error rate 2.991%

152 layers

Super-human precision

* Human-level performance: 5.1%

University of
Pittsburgh

# DL in speech recognition



According to Microsoft's speech group:

Using DL

Deep Learning in Speech Recognition

**BRIEF**

**Microsoft achieves 'human parity' in speech recognition system**

WER 5.1%

University of Pittsburgh

# Deep learning applications

### Image recognition



### Object detection



### Video game



### Image generation



### AlphaGO



### AlphaFold



### Medical diagnosis



### ChatGPT

# Categories of DL models

▸ **Supervised**

– Deep (Dense) Neural Networks (DNN)

– Convolutional Neural Networks (CNN)

– Recurrent (RNN); LSTM, transformers

▸ **Unsupervised**

– Auto-encoder (AE)

– Variational AE (VAE)

– Generative Adversarial Networks (GAN)

– Diffusion models

University of Pittsburgh

# Supervised learning (classification)

## Labeled training data (lots of)

*Training*



Cat

University of Pittsburgh

# Unsupervised DL models (Autoencoder, generative adversarial networks(GAN))

## Goal: Learn the codes (latent representation) of data

# Visualizing the code

# Reconstruct

# Generative Adversarial Networks (GAN)

**Goal:** Train a generator (decoder) G that learns from data to generate samples from data distribution

Training

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))].$$



University of Pittsburgh

# GAN Arts





**Artificial intelligence takes on song-composing duties in Eurovision-inspired contest**

BY RODRIGO PÉREZ ORTEGA
| APR. 24, 2020

# Deepfake

You can now watch the video.



BREAKING NOW
**DONALD TRUMP JOINS 2016 RACE FOR PRESIDENT** CNN
Donald Trump | (R) Presidential Candidate
S&P ▲ 11.86
SITUATION ROOM

https://d285xazlytdv8t.cloudfront.net/output.mp4

# Drug design

# Processes of Building Deep Learning

- **Observation/Data:** $D = \{y, x\}$
- ~~Feature extraction and selection~~
- Modeling
  - *Goal*: Model $D$
  - *Task: Define f*: $\qquad y = F(x; w)$

    DL

    (hyper) parameters

- Training
  - *Goal*: Infer $w$ (including hyper-parameters)
  - *Task 1*: *Loss function* - $L(D, w)$    Cross entropy; MSE
  - *Task 2*: *optimization*
    - $\hat{w} = argmin_w \, L(D, w)$    Stochastic gradient descent
- Performance assessment
  - accuracy, AUC, PR, mAP, …

University of Pittsburgh

# Classification

# Curve fitting (regression)



University of Pittsburgh

# Processes of Building Deep Learning

▸ **Observation/Data:** $D = \{y, x\}$

▸ ~~Feature extraction and selection~~

▸ **Modeling**

- *Goal*: Model $D$
- *Task: Define f :*      $y = F(x; w)$

DL

(hyper) parameters

▸ **Training**

- *Goal*: Infer $w$ (including hyper-parameters)
- *Task 1*: *Loss function* - $L(D, w)$     Cross entropy; MSE
- *Task 2*: *optimization*
  - ▸ $\hat{w} = argmin_w \, L(D, w)$     Stochastic gradient descent

▸ **Performance assessment**

- accuracy, AUC, PR, mAP, …

University of Pittsburgh

# Single neuron

F



▶ **Pre-activation**

| Input: | 0, | 8, | 15, | 22 |
|--------|----|----|-----|----|

| Output: | 32, | 46.4, | 59, | 71.6 |
|---------|-----|-------|-----|------|

➡ **a = x*1.8 + 32**

▶ **Activation**

Activation function (Sigmoid)

**F = g(a) = g(x*1.8 + 32)**



$$g(a) = \mathrm{sigm}(a) = \frac{1}{1+\exp(-a)}$$

University of Pittsburgh

# Activation Functions

## Activation function



**Sigmoid function:**
$$y = \frac{1}{1+e^{-y}}$$

**ReLU function:**
$$y = \begin{cases} y & y => 0 \\ 0 & y < 0 \end{cases}$$

(a) Step function or threshold function
(b) Sigmoid function y = 1/(1+e-x); <span style="color:red">takes a real-valued input and squashes it to range between 0 and 1</span>
(c) ReLU function
(d) Softmax

**Classification Problems**



$$\frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

| Output layer | | Probabilities |
|---|---|---|
| 1.3 | | 0.02 |
| 5.1 | | 0.90 |
| 2.2 | | 0.05 |
| 0.7 | | 0.01 |
| 1.1 | | 0.02 |

<span style="color:red">Softmax converts the input vector to a probabilistic domain. This is very important for us for the final output layer.</span>

University of Pittsburgh

# Single hidden layer NN or fully connected layers

– Hidden layer pre-activation:

$$a_i = \boldsymbol{w}_i^{(1)T} \boldsymbol{x} + b_i^{(1)}$$

or

$$\boldsymbol{a}^{(1)} = \boldsymbol{W}^{(1)} \boldsymbol{x} + \boldsymbol{b}^{(1)}$$

– Hidden layer activation:

$$\boldsymbol{h}^{(1)} = \boldsymbol{g}(\boldsymbol{a}^{(1)})$$

– Output layer pre-activation:

$$a^{(2)} = \boldsymbol{w}^{(2)T} \boldsymbol{h}^{(1)} + b^{(2)}$$

– Output layer activation:

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = o(a^{(2)})$$

O( ): output activation function

# Deep Neural Networks (DNN)

- Could have $L$ hidden layers:

  ‣ layer pre-activation for $k>0$ $(\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x})$

  $$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

  ‣ hidden layer activation ($k$ from 1 to $L$):

  $$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

  ‣ output layer activation ($k=L+1$):

  $$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$

# Processes of Building Deep Learning

- ‣ **Observation/Data:** $D = \{y, x\}$
- ‣ ~~Feature extraction and selection~~
- ‣ **Modeling**
  - • *Goal*: Model $D$
  - • *Task: Define f*: $\quad y = F(x; w)$    DL
    - (hyper) parameters

- ‣ **Training**
  - • *Goal*: Estimate **w** (including hyper-parameters)
  - • *Task 1*: *Loss function* - $L(D, w)$    Cross entropy; MSE
  - • *Task 2*: *optimization*
    - ‣ $\widehat{w} = argmin_w \, L(D, w)$    Stochastic gradient descent
- ‣ **Performance assessment**
  - • accuracy, AUC, PR, mAP, …

University of Pittsburgh

# Goal of DL training

– Determine model weights (W, b), based on training data



W, b?

```
def function(C):




    return F
```

Training data

Input:   0,    8,   15,   22,     38

Output:  32,  46.4,  59,  71.6,  100.4

# Ingredients of DL Training

– Training data (labeled data)

– Loss function $L(D, \boldsymbol{w})$

– Optimizers

$$\widehat{\boldsymbol{w}} = argmin_{\boldsymbol{w}}\ L(D, \boldsymbol{w})$$

# Loss functions

| | Hidden Layer | Output Layer (prediction) | Actual Output (label) |
|---|---|---|---|
| | $h = f(x) = w_{it_1} * x + b_1$ | 12.1 ⟷ | 32 |
| | $h = f(x) = w_{it_2} * x + b_2$ | 19.9 | 32 |
| | ... | | |
| | $h = f(x) = w_{it_n} * x + b_n$ | 31.4 | 32 |

– Assesses how good an estimate of (W, b) is
– Popular loss functions

    – Cross entropy loss (Classification)
    – Mean square error (regression; real-valued output)

University of
Pittsburgh

# Optimization

$$\widehat{\boldsymbol{w}} = argmin_{\boldsymbol{w}}\ L(D, \boldsymbol{w})$$

– Solution (Optimizer): Stochastic gradient descent algorithm
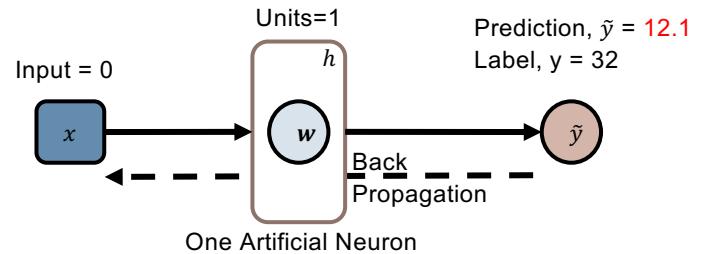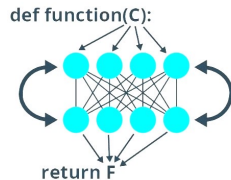– Efficient implementation: back-propagation

Optimizers

| Name | Abbreviation | Mean momentum | Std. momentum | Strengths |
|------|--------------|---------------|---------------|-----------|
| Stochastic Gradient Descent | SGD | FALSE | FALSE | Easy to understand |
| Nesterov Momentum | SGD w/ Nesterov | TRUE | FALSE | |
| Root Mean Square Propagation | RMSProp | TRUE | FALSE | Works well w/ text input |
| Adaptive Moment Estimation | Adam | TRUE | TRUE | Good default |

University of Pittsburgh

# What does it look like?

```
def function(C):
    F = C * 1.8 + 32
    return F
```

Input:   0,    8,    15,    22

Output:   32,   46.4,   59,   71.6
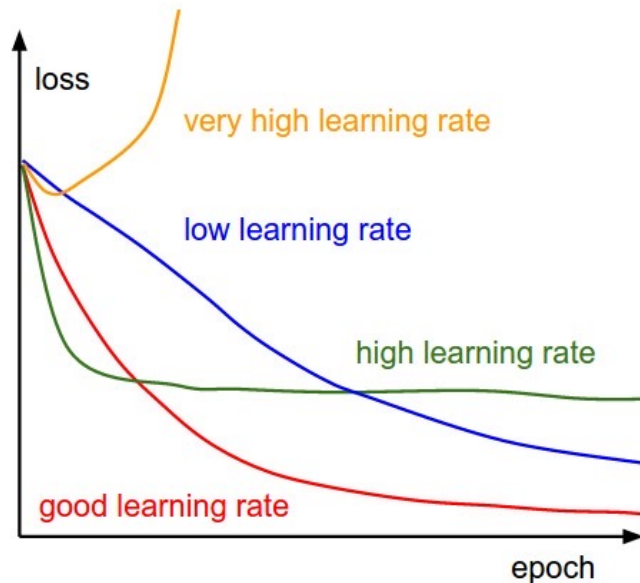
```
def function(C):
    ...
    return F
```

Units=1

Input = 0

$x$     $w$     $h$     $\tilde{y}$

Prediction, $\tilde{y}$ = 12.1
Label, y = 32

Back Propagation

One Artificial Neuron

| Iteration | Input Layer (input) | Hidden Layer | Output Layer (prediction) | Actual Output (label) | Compute Error (loss) | Loss function MSE |
|-----------|---------------------|--------------|---------------------------|-----------------------|----------------------|-------------------|
| 1 | 0 | $h = f(x) = w_{it_1} * x + b_1$ | 12.1 ⟷ | 32 | Loss(12, 32) -> very large | |
| 2 | 0 | $h = f(x) = w_{it_2} * x + b_2$ | 19.9 | 32 | Loss(19.9, 32) -> moderately large | |
| | | ... | | | | |
| n | 0 | $h = f(x) = w_{it_n} * x + b_n$ | 31.4 | 32 | Loss(31.4, 32) -> very small | |

University of Pittsburgh

# Important concepts

▸ Terminology

  ▸ **Batch size**: # of samples fed into an SGD step

  ▸ **Epoch**: # of steps that takes to use all training samples

  ▸ Learning rate

  ▸ Initial value

# Processes of Building Deep Learning

- ▸ Observation/Data: $D = \{y, x\}$
- ▸ ~~Feature extraction and selection~~
- ▸ Modeling
  - • *Goal*: Model $D$
  - • *Task: Define f* :     $y = F(x; w)$

    DL

    (hyper) parameters

- ▸ Training
  - • *Goal*: Estimate $w$ (including hyper-parameters)
  - • *Task 1*: *Loss function* - $L(D, w)$     Cross entropy; MSE
  - • *Task 2*: *optimization*
    - ▸ $\hat{w} = argmin_w\ L(D, w)$     Stochastic gradient descent
- ▸ Performance assessment
  - • accuracy, AUC, PR, mAP, …

# Evaluating classification performance

▸ Errors

   ▸ **False Positive**: Incorrectly labeled as relevant

   ▸ **False Negative**: Incorrectly labeled as not relevant

Prediction:

Image:



**True** Positive     **True** Negative     False Negative     False Positive
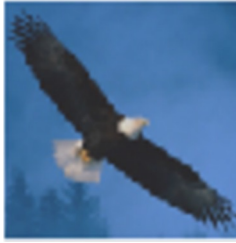
University of
**Pittsburgh**

# Evaluating classification performance

- Types of detection outcomes

| Decision | $\mathcal{H}_0$ (no signal) | $\mathcal{H}_1$ (having signal) |
|---|---|---|
| $\mathcal{H}_0$ | True negative | Miss, Type II, or F. negative |
| $\mathcal{H}_1$ | FA, Type I, or F. positive | Detection, True postive |

- Types of probabilities

| Decision | $\mathcal{H}_0$ (no signal) | $\mathcal{H}_1$ (having signal) |
|---|---|---|
| $\mathcal{H}_0$ | Specificity | $P_{Miss}, \beta$, False negative rate |
| $\mathcal{H}_1$ | $P_{FA}, \alpha$, False positive rate | $P_D, 1 - \beta$, power, sensitivity |

# Specificity and Sensitivity

▸ Specificity
  ▸ True negative probability;
  ▸ 1-False positive probability;
  ▸ Percentage of negative examples that are correctly labeled
  ▸ Specificity= (# true negatives) / (# negatives)

▸ Sensitivity (Recall)
  ▸ True positive probability;
  ▸ Percentage of positive examples that are correctly labeled
  ▸ Recall = (# true positives) / (# positives)

# Accuracy and precision

▸ Precision
  - ▸ Percentage of positive labels that are correct
  - ▸ Precision = (# true positives) / (# true positives + # false positives)

▸ Accuracy
  - ▸ Percentage of correct labels
  - ▸ Accuracy = (# true positives + # true negatives) / (# of samples)
  - ▸ Accuracy = 1 − P(error)

# Example

**Prediction:**



**Image:**

| True Positive | True Negative | False Negative | False Positive | False Positive | True Positive |

Specificity = 1/3      Recall = 2/3

Precision = 2/4      Accuracy = 3/6

University of Pittsburgh

# ROC curve and Area Under the Curve (AUC)
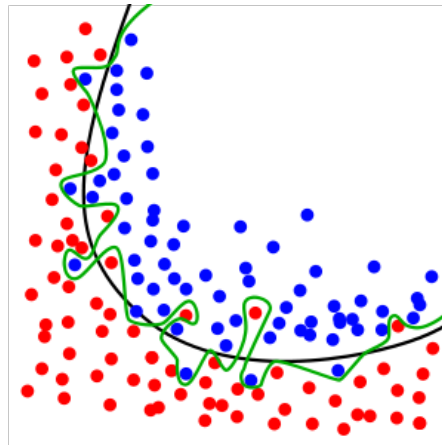
University of
Pittsburgh

# When to use which measure?

- No preferred labels and proportion of labels is unknown
  - ROC
- No preferred labels and proportion of labels is known
  - Accuracy
- Have a preferred label and proportion of the preferred label is small.
  - Precision vs. recall

University of
Pittsburgh

# Measuring Success/Failure for Classification
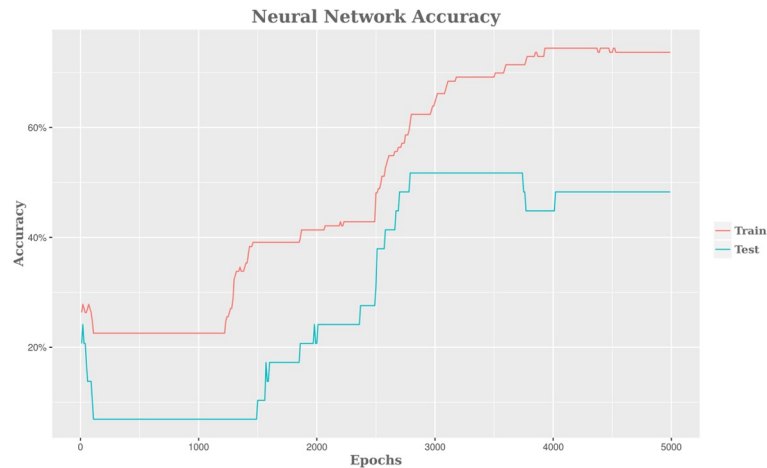
▸ Can we evaluate classification performance using training data?

▸ No, because these could be a classifier that can produce 0 error on training data. This is called overfitting.

▸ Overfitting
   ▸ Model performs well on training data but poorly on test data

# Use test data to measure success

▶ Training Data
  ▶ data used to learn a model
▶ Test Data
  ▶ data used to assess the accuracy of model

▶ What to do when you only have training data?

# Cross Validation

▸ **To avoid overfitting**

- train on part of available data, and test on rest
    - if dataset large (say, in 1000's), can simply set aside $\approx 1000$ random examples as test
    - otherwise, use 10-fold cross validation
        - break dataset randomly into 10 parts
        - in turn, use each block as a test set, training on other 9 blocks



training          testing

University of Pittsburgh