# Introduction to Neural Networks
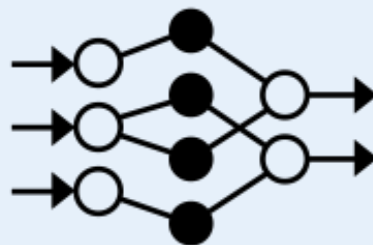
Qiyang Hu

UCLA Office of Advanced Research Computing

Oct 28th, 2022

# About this talk

- An introduction, an overview
  - The intuitive explanations on basic concepts
  - The advanced technical developments

- The outline
  - Machine learning and Deep learning
  - Neural network modeling in a general ML/DL workflow

- My DL talks in this and next quarters
  - Introduction to NN (today)
  - Learning PyTorch (next Wednesday)
  - Deep learning, the GBU (next Friday)
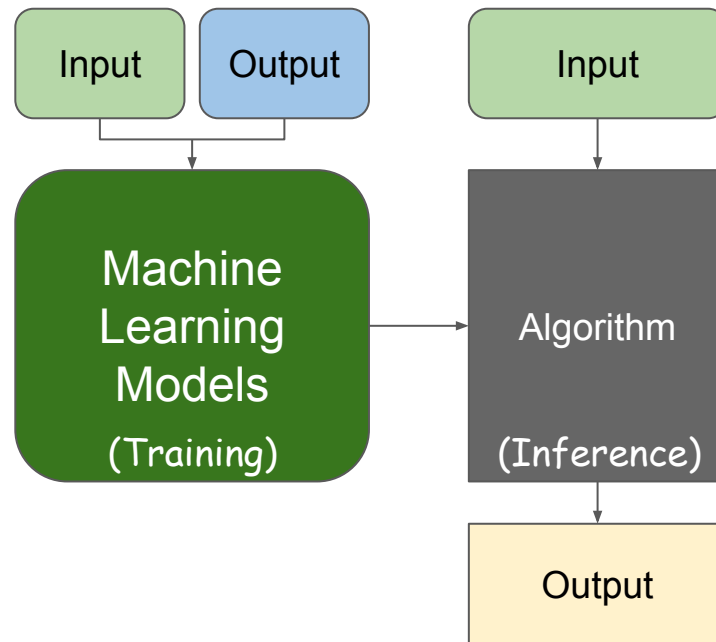  - Special NN topics, (conv, gans, transformer, lstm?) (next quarter)

Qiyang Hu
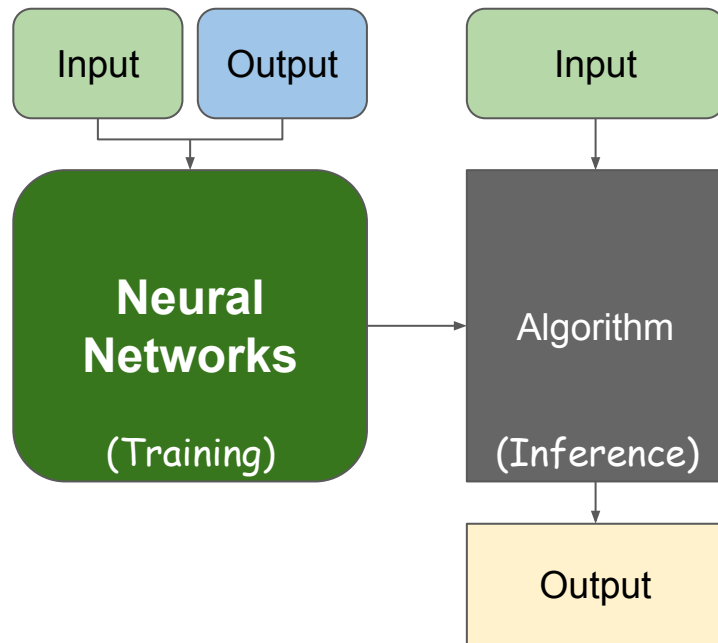
# What is Machine Learning?



Traditional Programming

Input → Known Algorithm → Output

Machine Learning

Input, Output → Machine Learning Models (Training) → Algorithm (Inference) ← Input → Output

Qiyang Hu

# What is   Deep   Learning?

**Traditional Programming**

**Deep Learning**

| Input |
|---|

| Known Algorithm |
|---|

| Output |
|---|

| Input | Output |
|---|---|

| Input |
|---|

**Neural Networks**

(Training)

Algorithm

(Inference)

| Output |
|---|

# Simplified workflow for a deep learning project



Step 1         Step 2    Step 3    Step 4    Step 5

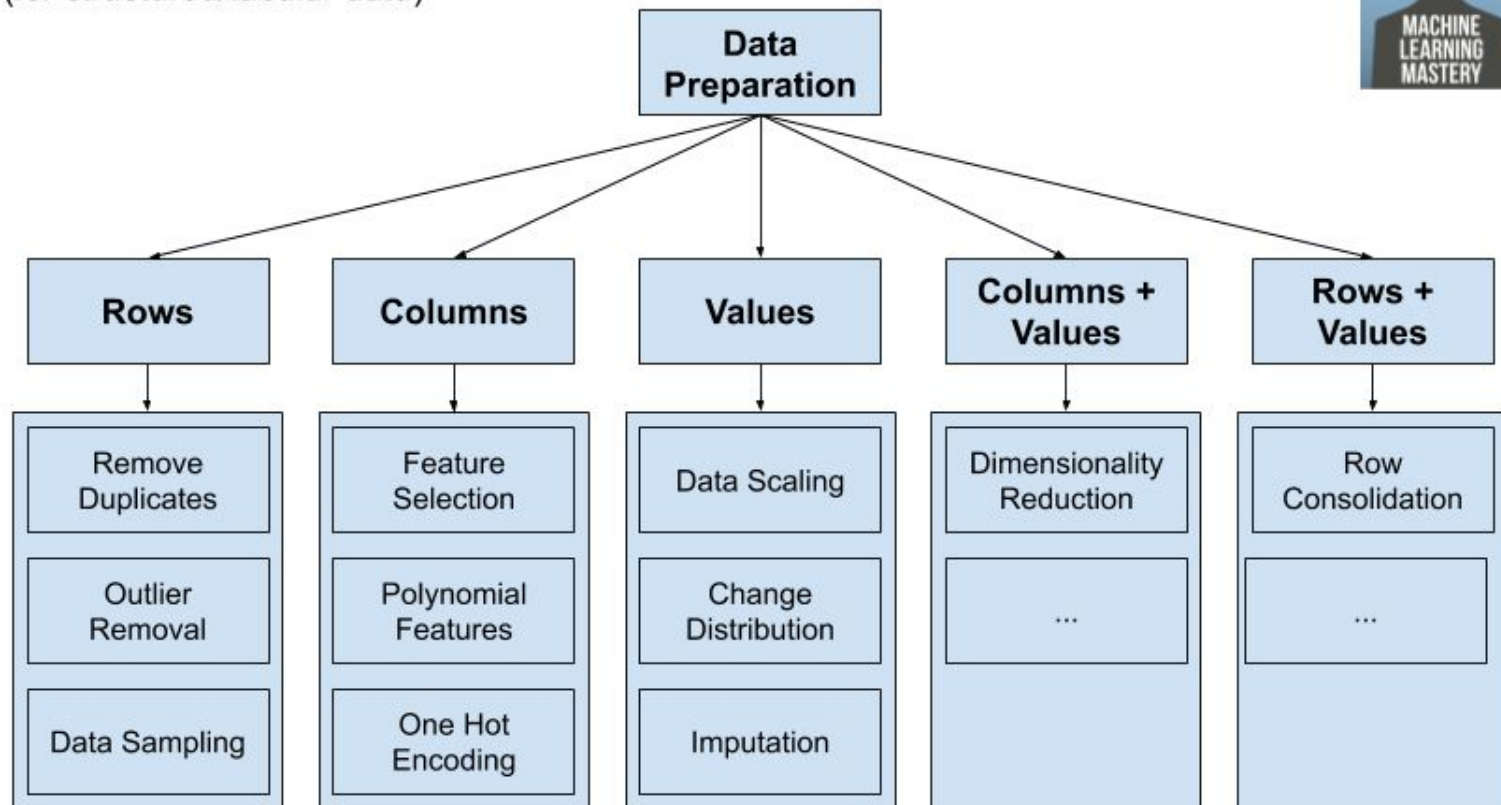Qiyang Hu

# Step 1. Data Preparation and Processing

- The most time-consuming but the most *creative* job
  - Take ~66% time
  - Require experience
  - May need domain expertise

- Determines the upper limit for the goodness of DL
  - Models/Algorithms: just approach the upper limit

What Data Scientists spend the most time doing

Deploying Models
11.0%

Data Loading
19.0%

Model Training and
12.0%

Model Selection
11.0%

Data Cleansing
26.0%

Data Visualization
21.0%

Anaconda's State of Data Science Report, 2020 ([Source](#))

Qiyang Hu

# Data Preparation Framework
(*for structured/tabular data*)



| Rows | Columns | Values | Columns + Values | Rows + Values |
|---|---|---|---|---|
| Remove Duplicates | Feature Selection | Data Scaling | Dimensionality Reduction | Row Consolidation |
| Outlier Removal | Polynomial Features | Change Distribution | ... | ... |
| Data Sampling | One Hot Encoding | Imputation | | |

Data Preparation

Qiyang Hu

# More data-prep tasks might be needed

- Image Data Processing
  - Pixel scaling
  - Train-Time Augmentation
  - Test-Time Augmentation
  - Convolution and Flattening

- Data Tokenization
  - Breaking the sequence data into units
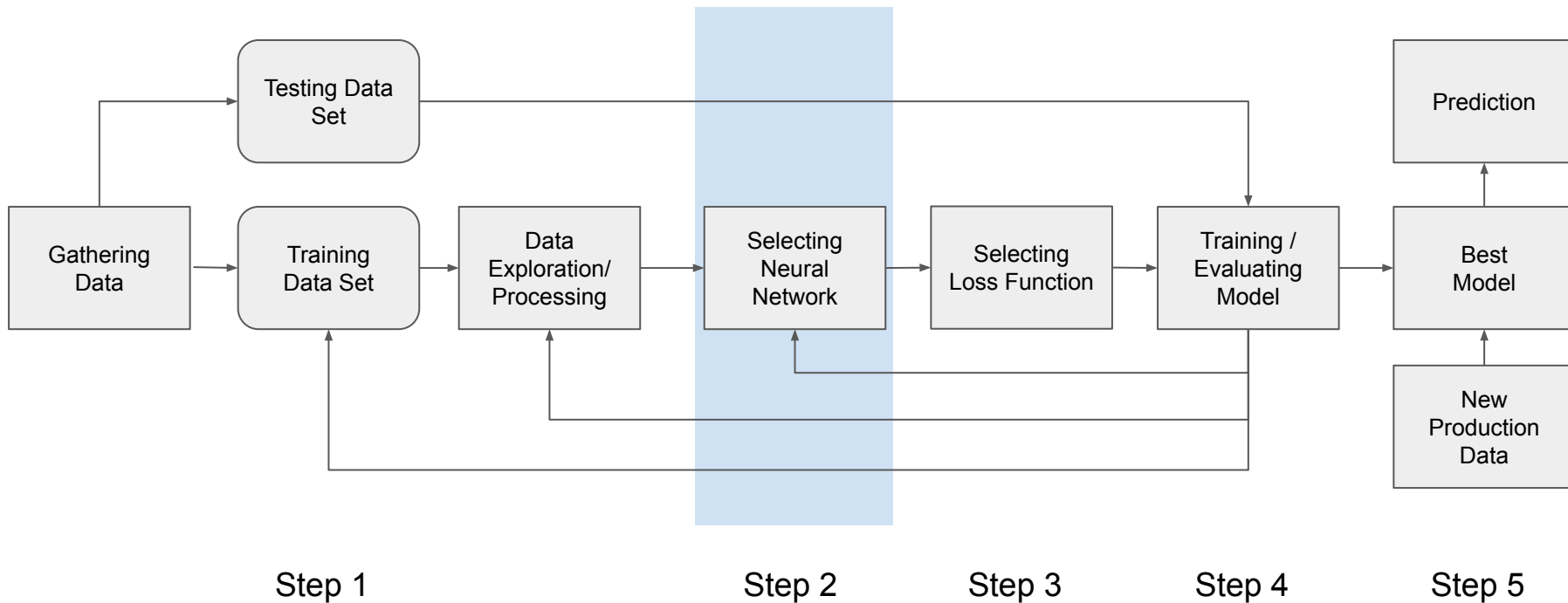  - Mapping units to vectors
  - Aligning & padding sequences

- Data Embedding
  - Map data to lower-dim vectors
    - Sparse to dense
    - Merging diverse data
    - Preserve relationship
  - Techniques
    - Std Dimensionality Reduction
    - Word2Vec
    - Be part of the model training
  - *Representation Learning*

  $$\text{Embedding Dims} \approx \sqrt[4]{\text{Possible Values}}$$
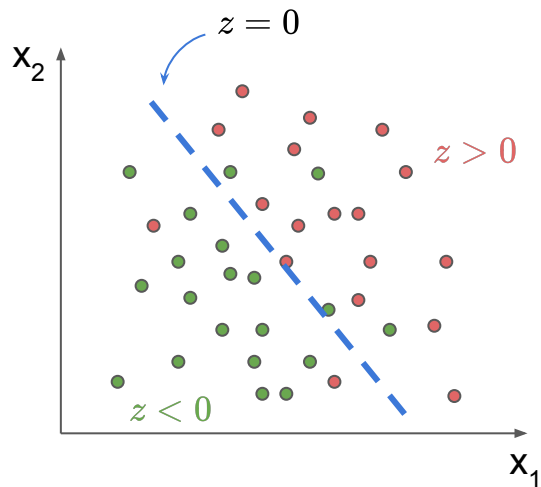
Qiyang Hu
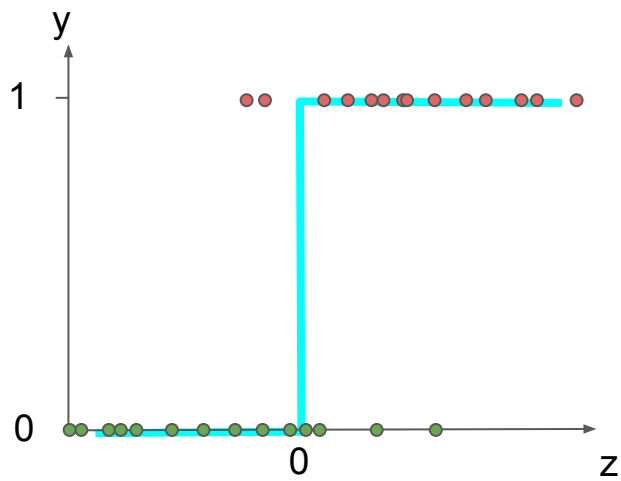
# Workflow for a deep learning project



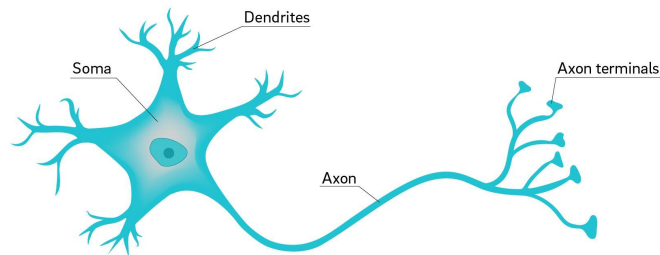Qiyang Hu

# What is Neural Network?

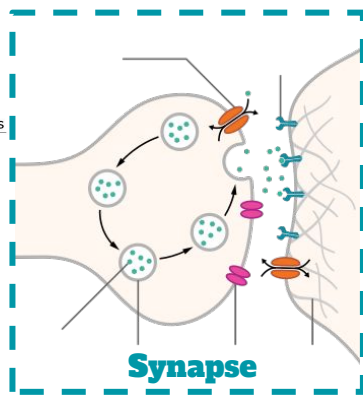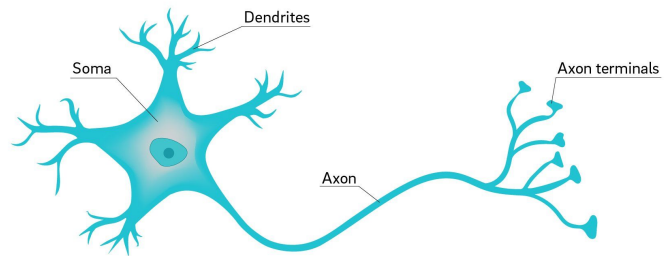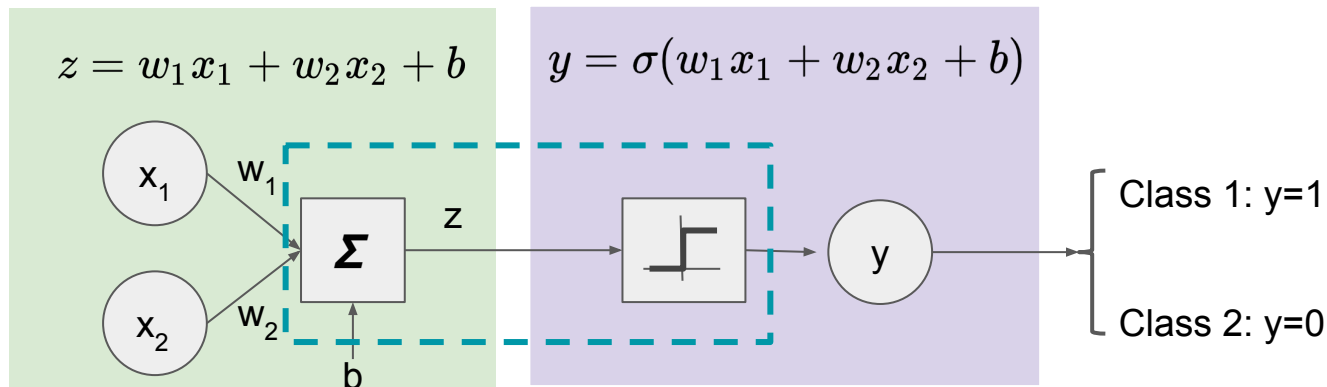- Recap for simple linear classification problem



$$z = w_1 x_1 + w_2 x_2 + b$$
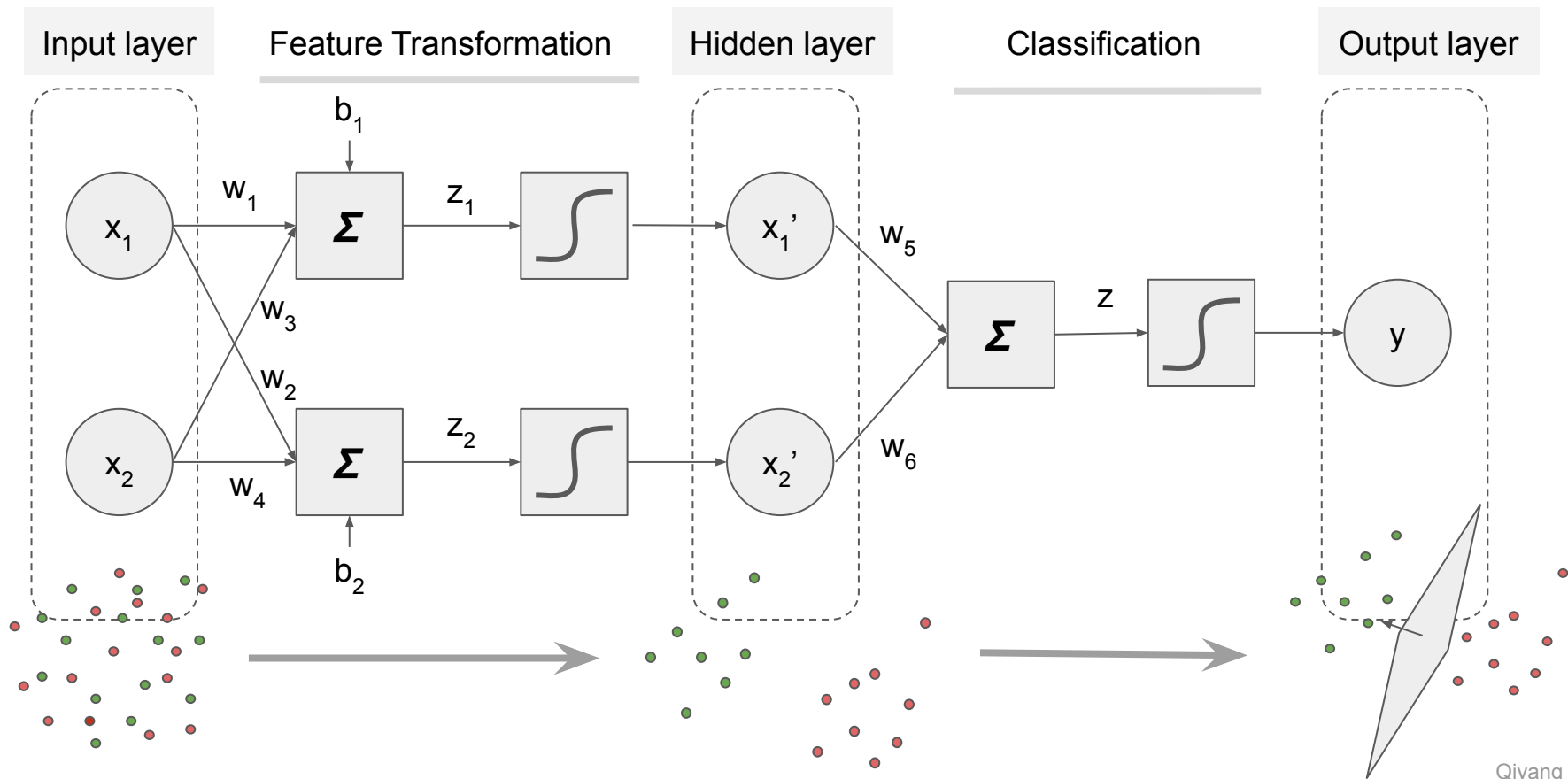
$$y = \sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \end{cases}$$

# Artificial Neuron and Biological Neuron

$$z = w_1x_1 + w_2x_2 + b$$

$$y = \sigma(w_1x_1 + w_2x_2 + b)$$

**McCulloch-Pitts (MCP) neuron model**

$x_1$   $w_1$

$x_2$   $w_2$

$\Sigma$   z

b

y

Class 1: y=1

Class 2: y=0

Dendrites

Soma

Axon terminals

Axon

Synapse

Dendrites

Soma

Axon terminals

Axon

Qiyang Hu

# Neural Networks ~ piling/stacking logistic-regression classifiers



Qiyang Hu

# *Deep* neural networks

- [LeNet-5](#) (1998)



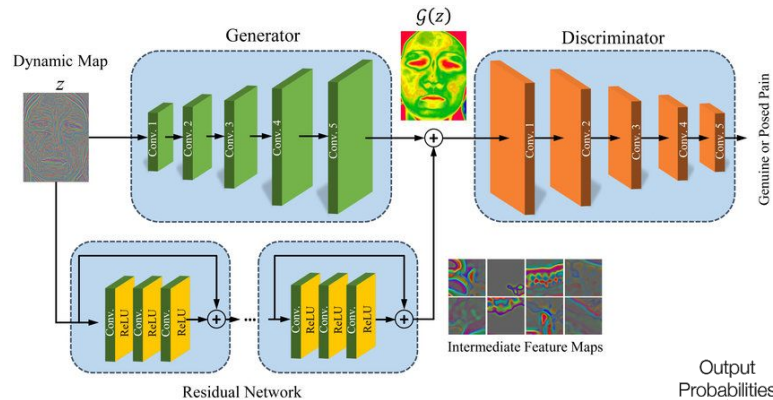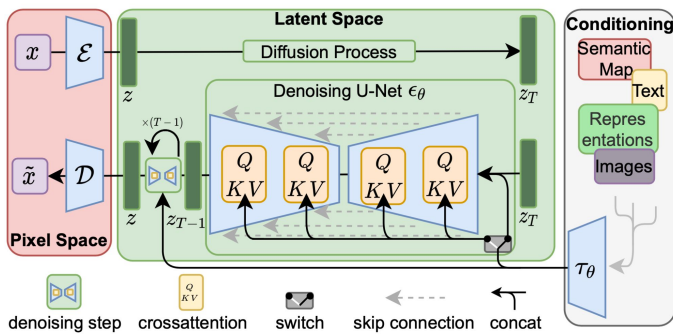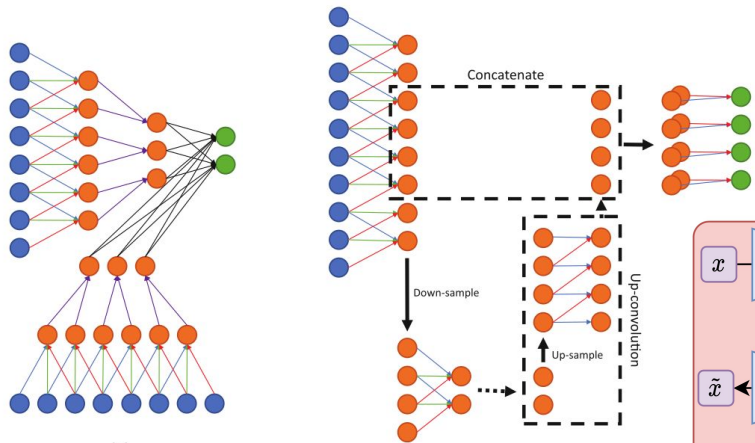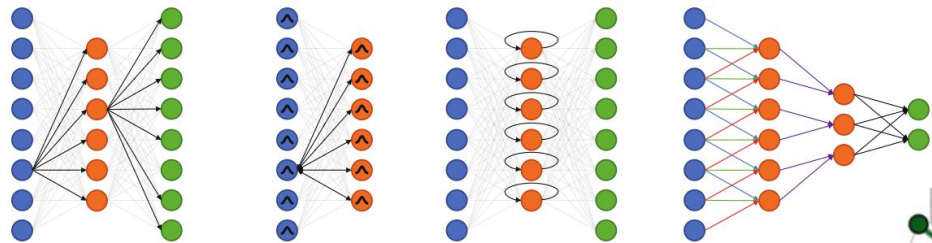| Year | CNN | Developed by | Place | Top-5 error rate | No. of parameters |
|------|-----|--------------|-------|------------------|-------------------|
| 1998 | LeNet(8) | Yann LeCun et al | | | 60 thousand |
| 2012 | AlexNet(7) | Alex Krizhevsky, Geoffrey Hinton, Ilya Sutskever | 1st | 15.3% | 60 million |
| 2013 | ZFNet() | Matthew Zeiler and Rob Fergus | 1st | 14.8% | |
| 2014 | GoogLeNet(19) | Google | 1st | 6.67% | 4 million |
| 2014 | VGG Net(16) | Simonyan, Zisserman | 2nd | 7.3% | 138 million |
| 2015 | ResNet(152) | Kaiming He | 1st | 3.6% | |

# Why deep?

- Shallow network can fit any function
  - Has less number of hidden layers
  - Has to be really "fat"

- Deep network is more efficient.
  - Exponentially fewer parameters (2017)
  - It can extract/build better features

# Deep neural *networks*



Qiyang Hu

# A simpler classification of neural network types

- Feed forward neural networks (No cycle in node connections)
  - Fully connected network
  - Convolutional networks (CNNs)



Input Layer  Hidden Layer  Output Layer

- Recurrent networks (w/ directed cycle in node connections)
  - Fully recurrent NN
  - Recursive NN
  - Long short-term memory (LSTM)
  - Hopfield network (w/o hidden nodes)



Input Layer  Hidden Layer  Output Layer

- Symmetric networks (no directions in node connections)
  - Boltzmann Machines
    - RBM, DBM, SOM



Visible Units  Hidden Units

# Activation Function

- **Sigmoid function:** $\sigma(z) = \dfrac{1}{1 + \exp(-z)}$

- **tanh function:** $tanh(z) = \dfrac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

- **Rectified linear unit (ReLU)**
  - Softplus  $\quad f(x) = x^+ = \max(0, x)$
  - Leaky ReLU
  - Exponential LU (ELUs)
  - GELU
  - Dynamic ReLU

- **Softmax function:** $y_i = \dfrac{e^{z^{(i)}}}{\sum_{j=0}^{K} e^{z^{(j)}}}$

- **Maxout Network:**
  - *Learnable* activation function

# How to choose activation functions?

For <u>hidden</u> layers



For <u>output</u> layers



From Machine Learning Mastery Blog Post

# Workflow for a deep learning project



Qiyang Hu

# How to measure the performance of the model?

- General name: objective function

- Measure the misfit of the model as a function of parameters
  - Criterion is to *minimize* the error functions
  - Loss Function, Cost Function: a penalty on difference between predictions and labels

- Evaluate the probability of *generating* training set
  - Criterion is to *maximize* the distribution likelihood as a function of parameters
  - Maximum (log)-likelihood estimation: minimize the divergence of distributions

- Regression losses and classification losses

Qiyang Hu

# Loss functions

- Generative/Predictive:



  - Regression Loss
    - Mean Square Error / Quadratic Loss / L2 Loss:
    - Mean Absolute Error / L1 Loss:
    - Huber Loss
    - Quantile Loss
  - Cross-Entropy Loss and variations
    - Log Loss / Negative Log Likelihood
    - Weighted CE / Balanced CE / Focal Loss
    - Dice Loss / IOU Loss / Tversky Loss

$$L_{MSE} = \frac{1}{n} \sum_{i}^{n} (t_i - s_i)^2$$

$$L_{MAE} = \frac{1}{n} \sum_{i}^{n} |t_i - s_i|$$

$$L_{CE} = - \sum_{i}^{C} t_i \log(s_i)$$

- Contrastive:



  - Ranking Loss/Margin Loss/Contrastive Loss/Triplet Loss

Qiyang Hu

# Workflow for a deep learning project



Step 1           Step 2    Step 3    Step 4    Step 5

# Training a DNN is an optimization problem



$$\theta = [w_1, w_2, \ldots, b_1, b_2, \ldots]$$

$$L^{\langle n \rangle}(\theta)$$

$$\sum_{n=1}^{N} L^{\langle n \rangle}(\theta)$$

$$\Rightarrow C(\theta)$$

- We know how to compute *C(θ)*, analytically or numerically.
- Start from an arbitrary initialization of $\theta_0$, and get an initial $C_0(\theta)$
- Apply optimization algorithm to minimize *C(θ)*

# Neural Network's Optimization

- Gradient Descent (a 1st-order approach)   $\theta \longleftarrow \theta - \eta \nabla L(\theta)$
  - Most popular algorithm
    - Pros: simple and fast
    - Cons: sometimes hard to tune

# Gradient-Descent Optimizers

- Stochastic GD / Mini-Batch GD
- Adding momentum:
  - Classical Momentum (CM)
  - Nesterov's Accelerated Gradient (NAG)
- Adaptive learning rate:
  - AdaGrad, AdaDelta, ...
  - RMSprop
- Combining the two
  - **ADAM** (as **default** in many libs)
- Beyond Adam:
  - Lookahead (2019), RAdam (2019)
  - AdaBound/AmsBound (ICLR 2019)
  - Range (2019)
  - AdaBelief (NeurIPS 2020 Spotlight)

Gradient descent vs Momentum vs AdaGrad vs RMSProp vs Adam

(Source)

Qiyang Hu

# Higher Order Optimization Algorithms

- Newton-like methods (2nd-order methods)

$$\theta \longleftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

  - Prons: **fewer** iterations, fewer hyperparameters
  - Cons: much more **costly** in each iteration, more storing

  - DFP/Broyden/BFGS/L-BFGS: a quasi-newton one
    - Good for low dimensional models

  - Conjugate gradient (CG): between GD and Newton
    - moderately high dimensional models

- Natural gradient descent methods  $\nabla_\theta L(\theta) = F^{-1} \nabla_\theta L(\theta)$
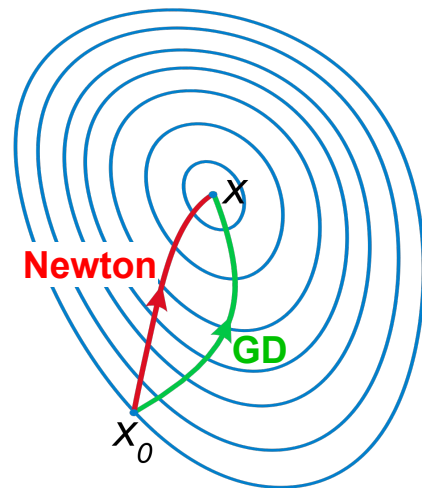
  - K-FAC (Martens and Grosse, 2015)
  - Shampoo (Gupta, et al., 2018)
  - K-BFGS (Goldfarb, et al., NeurIPS 2020)



Figure from [Wikipedia](Wikipedia)

Qiyang Hu

# Using Gradient Descent to train DNN



$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}^{\langle n \rangle}$$

$$\theta = [w_1, w_2, \ldots, b_1, b_2, \ldots]$$

Millions of parameters!

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix}^{\langle n \rangle}$$

$$L^{\langle n \rangle}(\theta)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \end{bmatrix}^{\langle n \rangle}$$

$$N$$

$$\sum_{n=1}^{N} L^{\langle n \rangle}(\theta)$$

$$\Rightarrow C(\theta)$$

$$\theta_0 \rightarrow \nabla C(\theta_0) \rightarrow \theta_1 \rightarrow \nabla C(\theta_1) \rightarrow \theta_2 \rightarrow \ldots$$

$$\theta_1 = \theta_0 - \eta \nabla C(\theta_0)$$

$$\theta_2 = \theta_1 - \eta \nabla C(\theta_1)$$

$$\vdots$$

$$\nabla C(\theta) = \sum_{n=1}^{N} \begin{bmatrix} \dfrac{\partial L^{\langle n \rangle}(\theta)}{\partial w_1} \\ \dfrac{\partial L^{\langle n \rangle}(\theta)}{\partial w_2} \\ \vdots \\ \dfrac{\partial L^{\langle n \rangle}(\theta)}{\partial b_1} \\ \vdots \end{bmatrix}$$

How to compute the gradient vector with millions of elements **efficiently?**

Qiyang Hu

# Backpropagation: a game of chain rule



$$y = \sigma_L \left( w_L \cdot \sigma_{L-1} \left( \cdots w_2 \cdot \sigma_1 (w_1 \cdot x + b_1) + b_2 \right) + b_L \right)$$

$$\underbrace{w_1 \cdot x + b_1}_{z_1}$$

$$\underbrace{\phantom{w_1 \cdot x + b_1}}_{a_1}$$

$$\frac{\partial C(y(w) - \hat{y})}{\partial w} = \frac{\partial z}{\partial w} \frac{\partial C}{\partial z} = \frac{\partial z}{\partial w} \left[ \frac{\partial a}{\partial z} \boxed{\frac{\partial C}{\partial a}} \right] = \frac{\partial z}{\partial w} \left[ \sigma' \cdot \left( \frac{\partial z_{(+1)}}{\partial a} \frac{\partial C}{\partial z_{(+1)}} \right) \right]$$
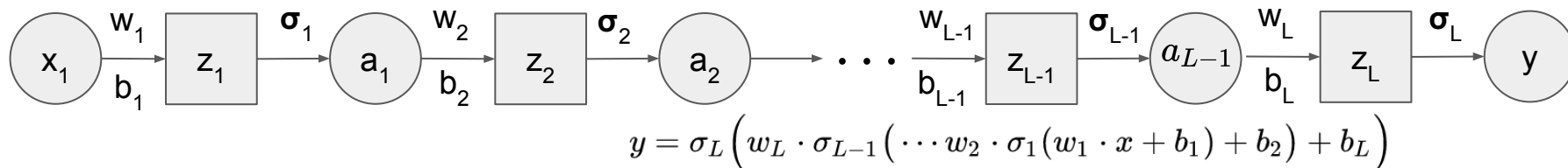
① Forward Pass

$$\frac{\partial z_1}{\partial w_1} = x_1 \longrightarrow \frac{\partial z_2}{\partial w_2} = a_1 \longrightarrow \cdots \longrightarrow \frac{\partial z_{L-1}}{\partial w_{L-1}} = a_{L-2} \longrightarrow \frac{\partial z_L}{\partial w_L} = a_{L-1}$$

② Backward Pass

$$\frac{\partial C}{\partial z_1} = \sigma'_1 \left[ w_2 \frac{\partial C}{\partial z_2} \right] \longleftarrow \cdots \longleftarrow \frac{\partial C}{\partial z_{L-1}} = \sigma'_{L-1} \left[ w_L \frac{\partial C}{\partial z_L} \right] \longleftarrow \frac{\partial C}{\partial z_L} = \sigma'_L \frac{\partial C}{\partial y} \longleftarrow \frac{\partial C}{\partial y}$$

# Gradient vanishing/exploding in DL training

- **Causes**



$$y = \sigma_L \left( w_L \cdot \sigma_{L-1} \left( \cdots w_2 \cdot \sigma_1 (w_1 \cdot x + b_1) + b_2 \right) + b_L \right)$$
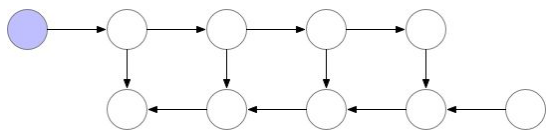
*After backprop* ←

  - Gradients in initial layers = Multiplication of Gradients at prior layers
  - Small variation around 1 results in vanishing/exploding
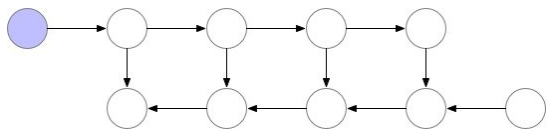
- **Techniques to resolve:**
  - General: adjusting learning rate, dropout, batch normalization, layer normalization
  - For gradient exploding: gradient clipping, weight regularization
  - For gradient vanishing: activation function, proper initialization parameters, LSTM, skip connections

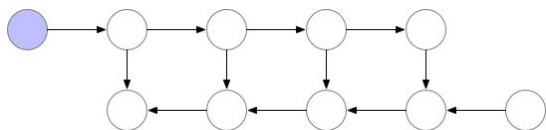# Backprop beyond the traditional neural networks

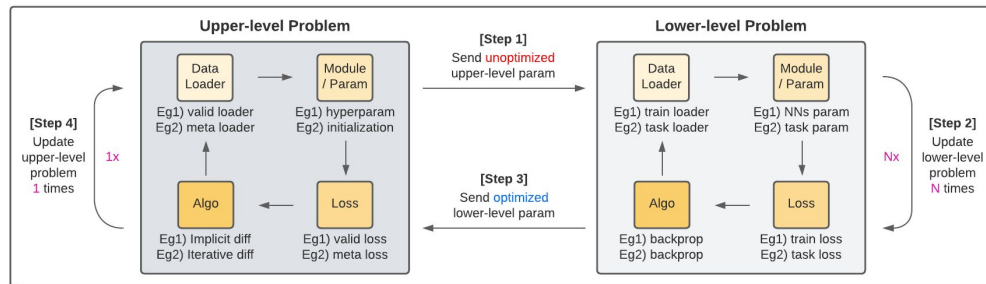- Gradient checkpointing ([source](#))



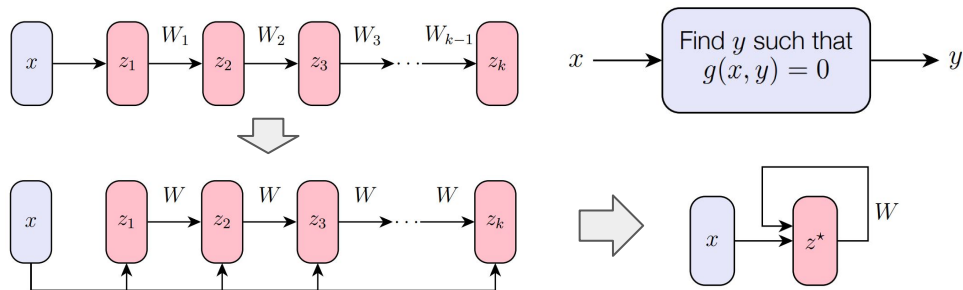Vanilla Backprop

Memory-poor Backprop

Checkpointed Backprop
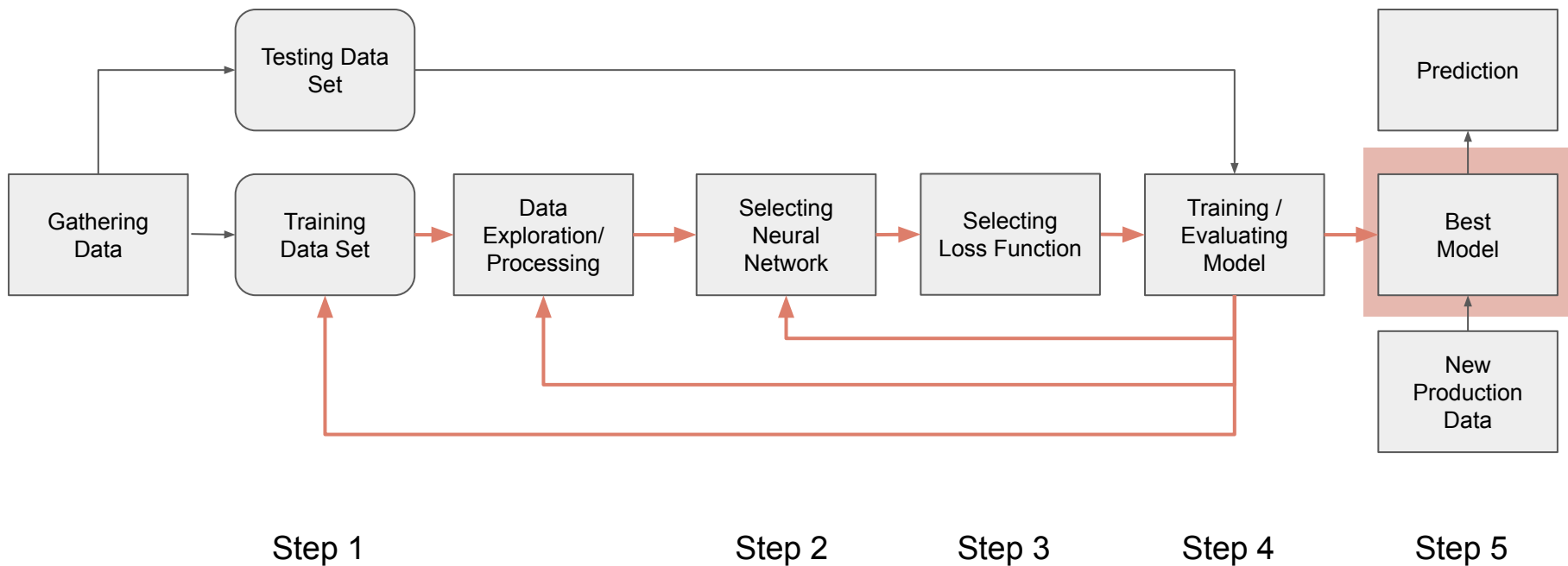
- Multi-level optimization



- Deep implicit layers



Qiyang Hu

# Backprop beyond deep learning

- A different way to calculate the differentiation of iterative math expressions
  - Not approximate, unlike numerical differentiation
  - Exact, like manual or symbolic differentiation, but with constant overhead

- Automatic differentiation (algorithmic differentiation)
  - Problems constructed by differentiable directed graphs (e.g. NN)
  - General functional blocks (FF, conv, recurrent blocks, etc)
  - Modularized optimization: differentiable optimizations in layer levels

- Differentiable physics
  - Physics problems represented by a sequence of differentiable operators
  - Differentiable programming
    - Enables classical numerical algorithms
    - Beyond simple chained transformations to include more complex control structures

# Workflow for a deep learning project



Testing Data Set

Prediction

Gathering Data

Training Data Set

Data Exploration/ Processing

Selecting Neural Network

Selecting Loss Function

Training / Evaluating Model

Best Model

New Production Data

Step 1          Step 2   Step 3   Step 4   Step 5

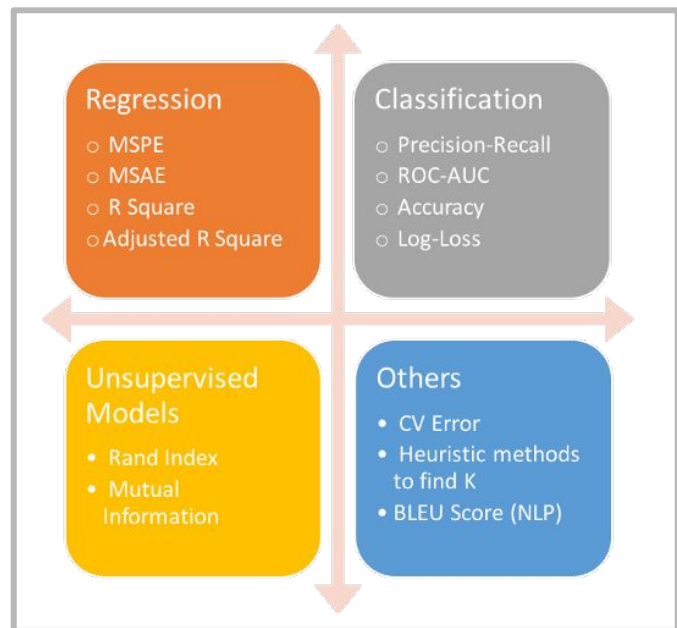Qiyang Hu

# Neural Network Evaluation

- Errors in regression (e.g. MSE):
  - From Model: features, algorithm ⇒ Bias
  - From Data: insufficient observations ⇒ Variance
  - From Noise

- Bias-Variance Trade-off

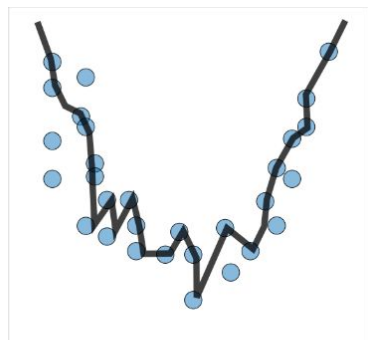$$\texttt{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$
$$\mathbb{E}[(\hat{\mu} - \mu)^2] = (\mathbb{E}[\hat{\mu} - \mu])^2 + \text{Var}[\hat{\mu} - \mu]$$





Regression
- MSPE
- MSAE
- R Square
- Adjusted R Square

Classification
- Precision-Recall
- ROC-AUC
- Accuracy
- Log-Loss

Unsupervised Models
- Rand Index
- Mutual Information

Others
- CV Error
- Heuristic methods to find K
- BLEU Score (NLP)

Qiyang Hu

# Underfitting and Overfitting

- Underfitting: model too simple:
  - Diagnose:
    - cannot even fit the training data
    - training error ~ testing error
  - Ignore the variance in training data
  - Higher prediction bias



- Overfitting: model too complex
  - Diagnose:
    - well-fit for training data
    - large error for testing data
  - Over-interpret training data
  - More deviation from new data



Qiyang Hu

# How to prevent

- Redesign the model
- Increase model's complexity
- Add more features as input
- Training longer

**Underfitting**

**Overfitting**

- Get more data
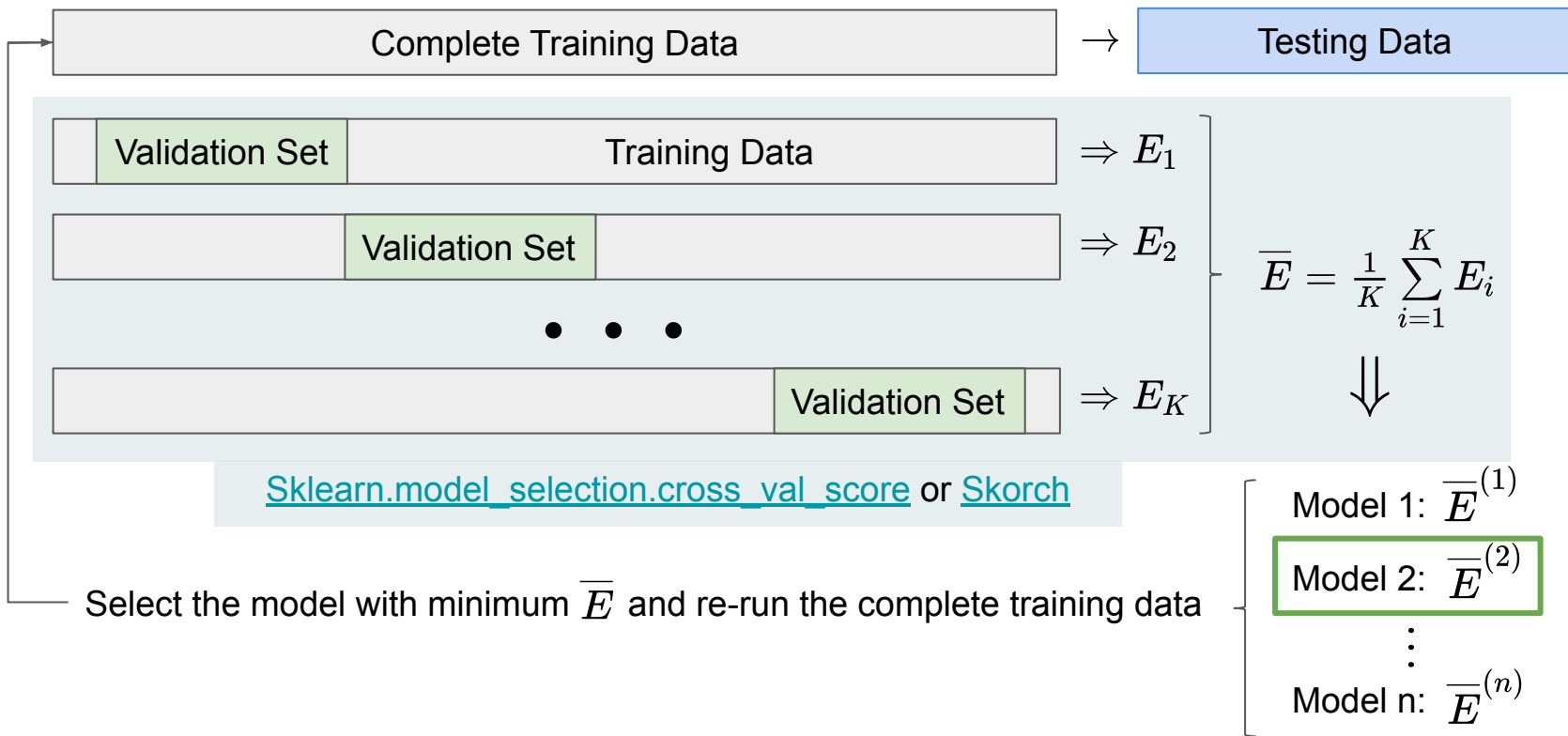  - collection or augmentation
- Reduce the model's complexity
- Regularization
  - Weight Regularization
  - Early stopping

Qiyang Hu

# Model Selection: K-fold Cross Validation



Complete Training Data $\rightarrow$ Testing Data

| Validation Set | Training Data | $\Rightarrow E_1$ |

| Validation Set | $\Rightarrow E_2$ |

$\bullet \quad \bullet \quad \bullet$

| Validation Set | $\Rightarrow E_K$ |

$$\overline{E} = \frac{1}{K} \sum_{i=1}^{K} E_i$$

$\Downarrow$

Sklearn.model_selection.cross_val_score or Skorch

Select the model with minimum $\overline{E}$ and re-run the complete training data

Model 1: $\overline{E}^{(1)}$

Model 2: $\overline{E}^{(2)}$

$\vdots$

Model n: $\overline{E}^{(n)}$

Qiyang Hu

# Errors/scores in practice



|  | | Public | Private |
|---|---|---|---|
| Training Set | Validation Set | Testing Set | Testing Set |

Error: $\quad E^{val} \quad < \quad E^{Pub} \quad < \quad E^{Pri}$

Score: $\quad S^{val} \quad > \quad S^{Pub} \quad > \quad S^{Pri}$

Qiyang Hu

# OARC Workshop Survey

http://bit.ly/3Dhp91H