



Tsinghua Deep Learning Summer School
2018/7/23

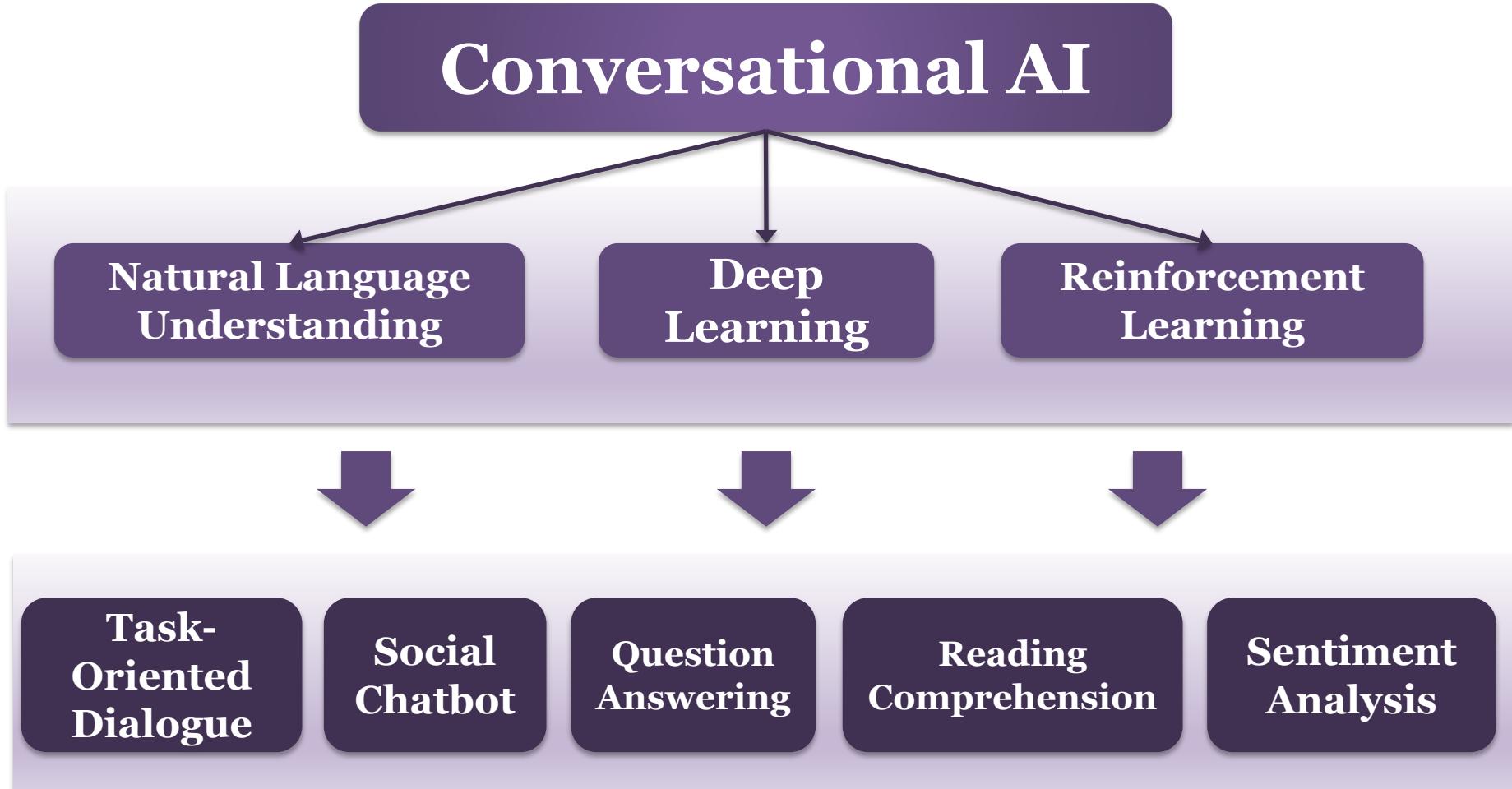
Introduction to Deep Learning

Minlie Huang

aihuang@tsinghua.edu.cn

<http://coai.cs.tsinghua.edu.cn/hml>

Research Picture of the CoAI Group



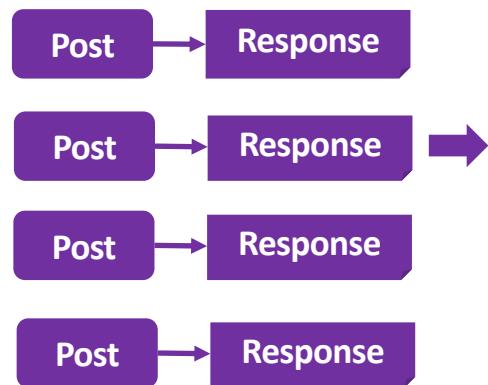
Major Research Directions

- **Sentiment & Emotion Understanding**
 - EMNLP 2016, ACL 2017, TOIS 2017, WWW 2018/2019
- **Language Generation**
 - ACL 2018, AAAI 2018/2019, IJCAI 2018/2019, WWW 2019
- **Dialog Systems**
 - ACL 2018, SIGIR 2018/2019, AAAI 2018/2019, IJCAI 2018/2019
- **Reinforcement Learning in NLP & Search**
 - WWW 2018/2019, AAAI 2018/2019, IJCAI 2018/2019

Emotional Chatting Machine

Perceiving and Expressing emotion by machine
Closer to human-level intelligence

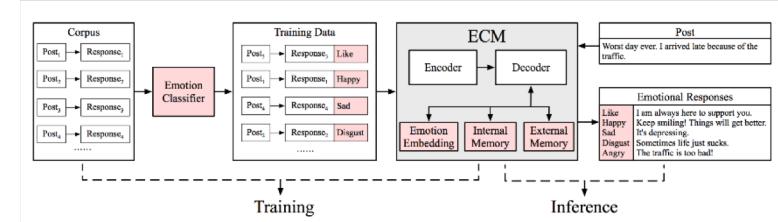
Social Interaction Data



Emotion Tagged data

今天领导前来视察工作！
Happy: 哈哈，领导视察工作辛苦了
Like: 这是一个很有特色的领导人
Sad: 我想知道，谁来帮我干活？
Angry: 这是要干什么？领导？

Emotional Chatting Machine

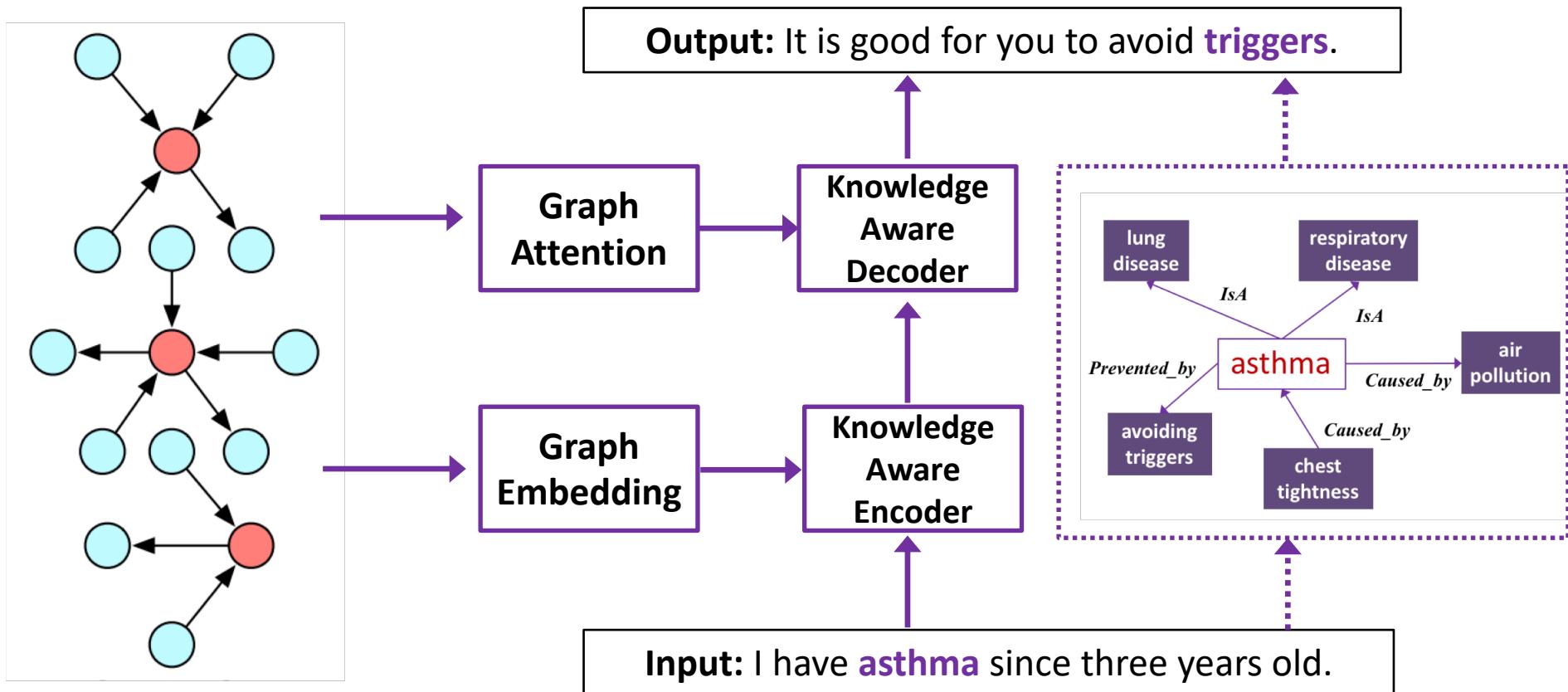


Our work was reported by **MIT Technology Review, the Guardian, Cankao News, Xinhua News Agency etc.**

Prof Björn Schuller: “an important step” towards personal assistants that could read the emotional undercurrent of a conversation and respond with something akin to empathy.

- Hao Zhou, Minlie Huang, Xiaoyan Zhu, Bing Liu. Emotional Chatting Machine: Emotional Conversation Generation with Internal and External Memory. **AAAI 2018**.

Commonsense-aware Chatbot



Commonsense Knowledge Aware Conversation Generation with Graph Attention. **IJCAI-ECAI 2018 distinguished paper** (3470 submissions; 710 accepted)

Outline

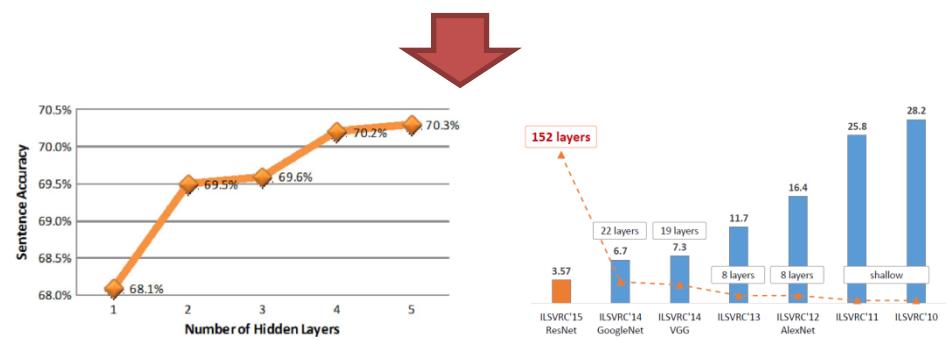
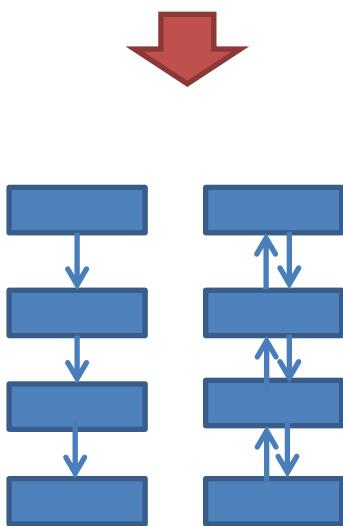
- Overview of deep learning
- Math basics
- Machine learning basics

Deep Learning

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.



What is deep learning and why is it so popular?

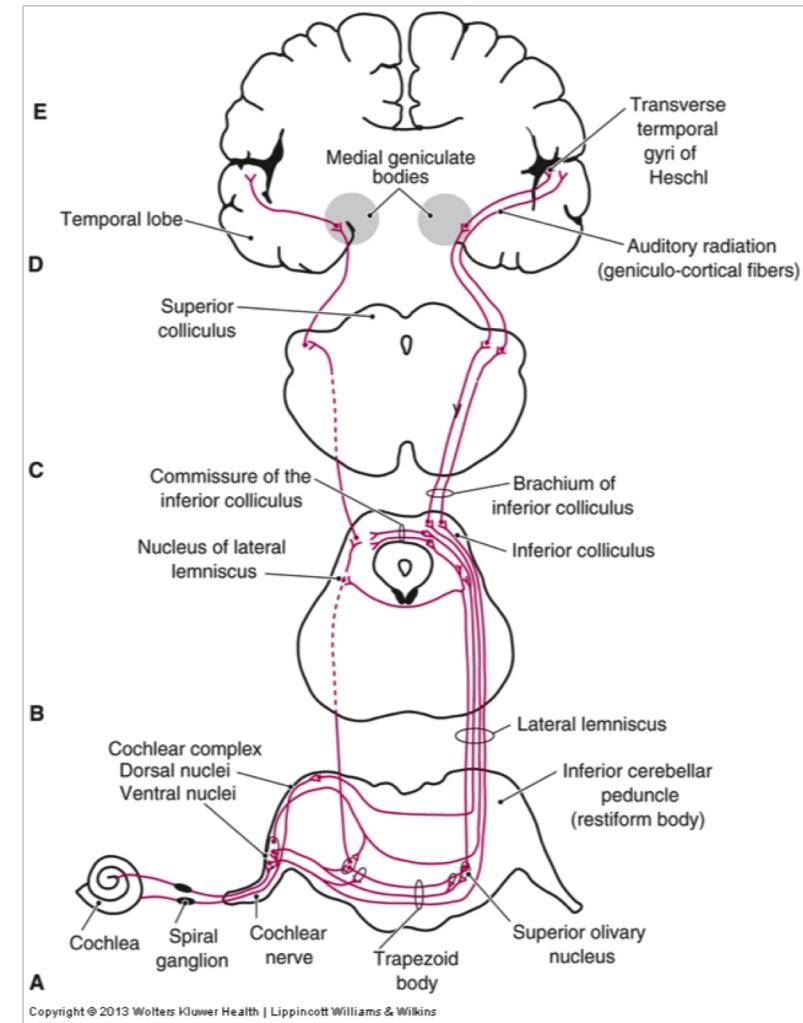
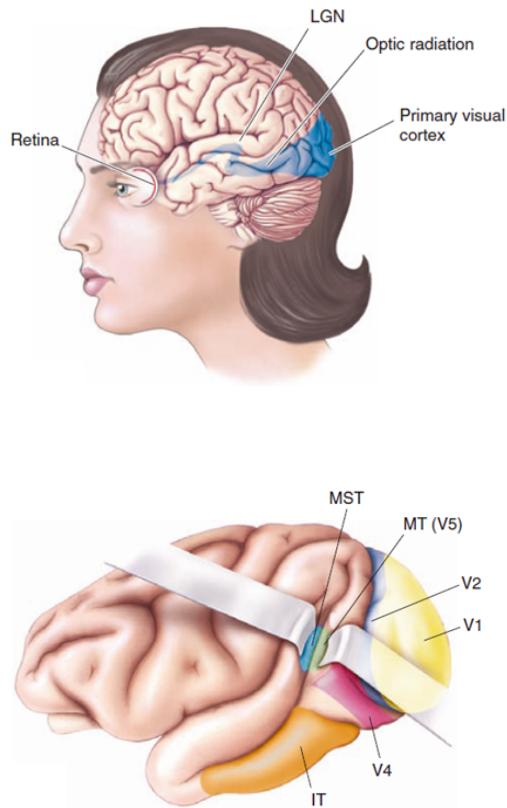


Why go deep?

- Data are often high-dimensional.
- There is a huge amount of **structure** in the data, but the structure is too complicated to be represented by a simple model.
- Insufficient depth can require more computational elements than architectures whose depth matches the task.
- Deep nets provide simpler but more descriptive model of many problems.

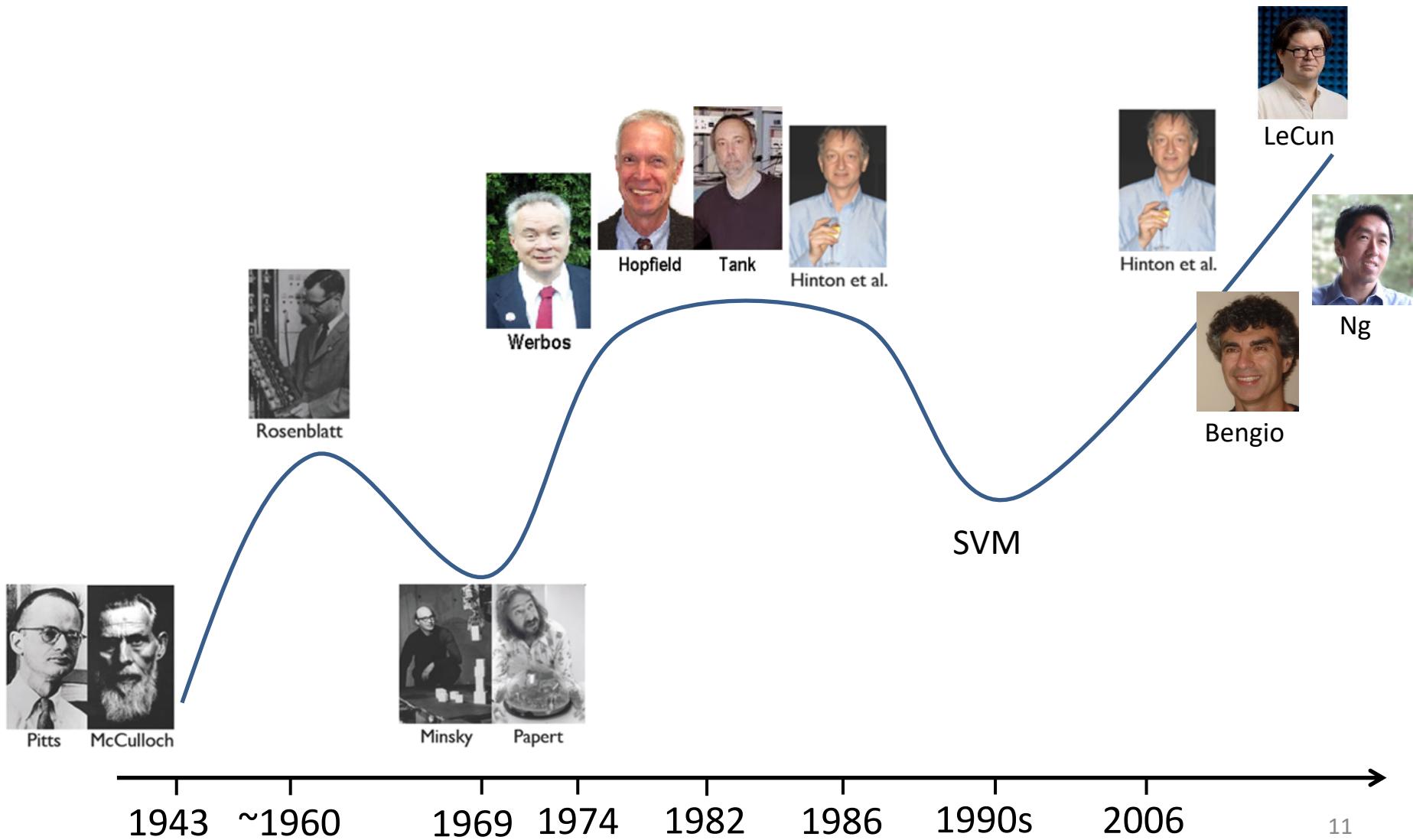
-By Geoffery Hinton

Hierarchical structures in the brain

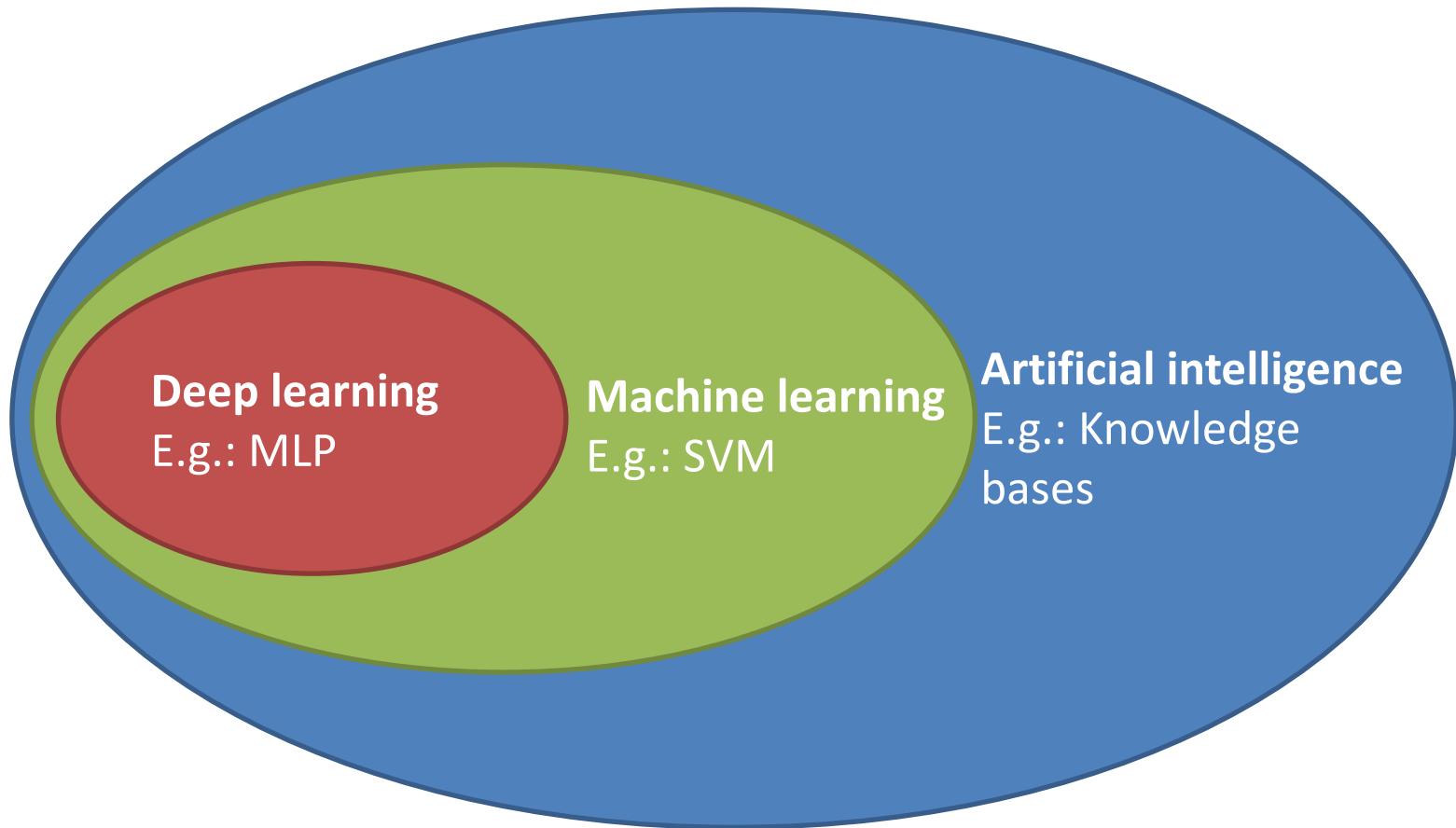


Copyright © 2013 Wolters Kluwer Health | Lippincott Williams & Wilkins

History of deep learning

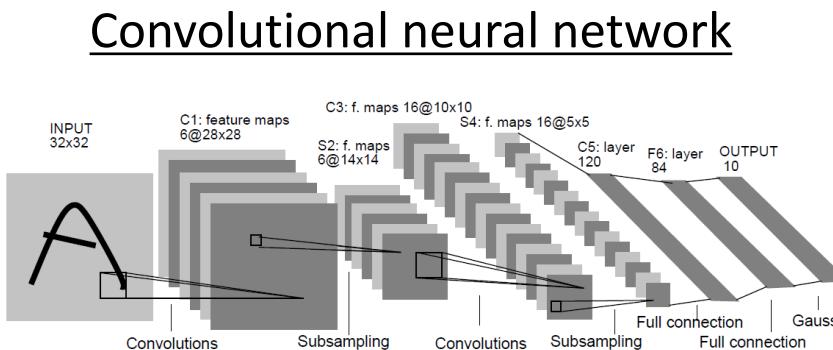
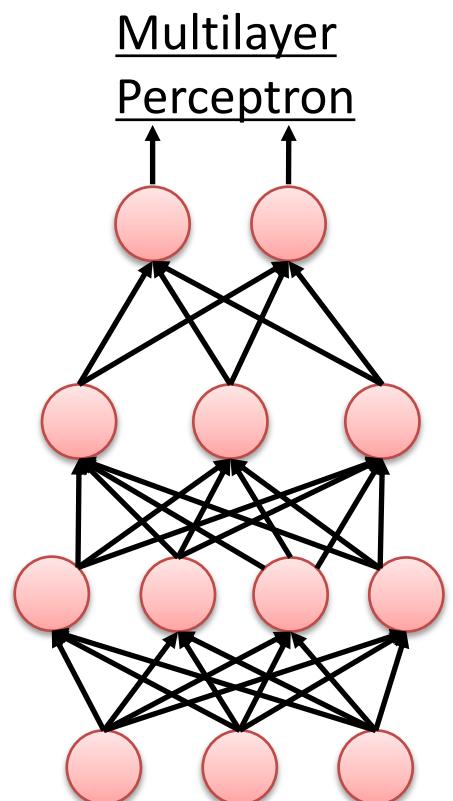


Deep learning, machine learning and AI

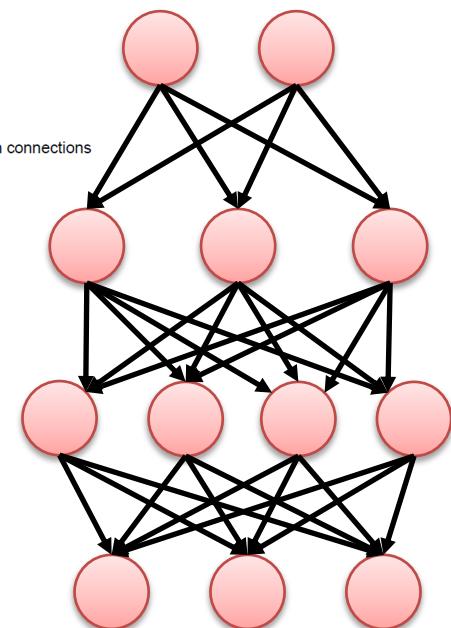


Contents to be covered in this summer school

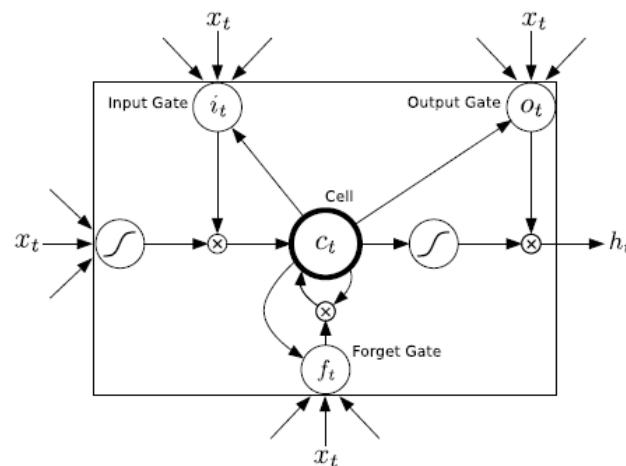
- Various deep learning models



Deep belief network



Recurrent neural network



Contents to be covered in this summer school

- Various applications
 - Image classification ([Xiaolin Hu & Wenguang Chen](#))
 - Social multimedia analysis ([Peng Cui](#))
 - Natural language processing ([Minlie Huang](#))
 - Knowledge graph ([Zhiyuan Liu](#))
 - Affective Computing ([Jia Jia](#))
 - Robot sensing ([Huaping Liu](#))
 - Deep Generative Model & Deep Reinforcement Learning ([Jun Zhu](#))

Outline

- Overview of deep learning
- **Math basics**
- Machine learning basics
- A brief intro to Python

Math basics

LINEAR ALGEBRA

Math Concepts

- Scalar
 - A single number, often denoted by a lower case letter without boldface, e.g., a, b, x

- Vector
 - An array of numbers, often denoted by a lowercase letter with boldface, e.g., $\mathbf{a}, \mathbf{b}, \mathbf{x}$

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}$$

- Matrix
 - A 2D array of numbers, often denoted by an uppercase letter with boldface, e.g., $\mathbf{A}, \mathbf{B}, \mathbf{X}$

$$\mathbf{A} = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix}$$

- Tensor
 - An n -D array of numbers, often denoted like this:
e.g., $\mathbf{A}, \mathbf{B}, \mathbf{X}$

$$\mathbf{A} = \left(\left(\begin{array}{cc} A_{1,1,1} & A_{1,2,1} \\ A_{2,1,1} & A_{2,2,1} \end{array} \right), \left(\begin{array}{cc} A_{1,1,2} & A_{1,2,2} \\ A_{2,1,2} & A_{2,2,2} \end{array} \right) \right)$$

Simple operations

- Matrix transpose: \mathbf{A}^\top

$$\mathbf{A} = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \quad \mathbf{A}^\top = \begin{pmatrix} A_{1,1} & A_{2,1} \\ A_{1,2} & A_{2,2} \end{pmatrix}$$

- A vector can be viewed as a special matrix

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \quad \mathbf{a}^\top = (a_1, a_2, a_3)$$

- Matrix product: if $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, then $\mathbf{C} = \mathbf{AB}$ with shape $m \times p$ and $C_{i,j} = \sum_k A_{i,k}B_{k,j}$
- Elementwise product (Hadamard product): $\mathbf{C} = \mathbf{A} \odot \mathbf{B}$ where the 3 matrices are of the same shape and $C_{i,j} = A_{i,j}B_{i,j}$

Simple operations

- Properties
 - $A(B + C) = AB + AC$
 - $A(BC) = (AB)C$
 - $AB \neq BA$ but $x^\top y = y^\top x$
 - $(AB)^\top = B^\top A^\top$
- Identity matrix: $I \in \mathbb{R}^{n \times n}$
 - For $x \in \mathbb{R}^n$, we have $Ix = x$
 - For $A \in \mathbb{R}^{n \times m}$, we have $IA = A$
- Matrix inverse of a square matrix A is a matrix denoted by A^{-1} such that $A^{-1}A = I$

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

Linear equations

- A system of linear equations

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{x} \in \mathbb{R}^n$, or equivalently

$$A_{1,1}x_1 + A_{1,2}x_2 + \cdots + A_{1,n}x_n = b_1$$

$$A_{2,1}x_1 + A_{2,2}x_2 + \cdots + A_{2,n}x_n = b_2$$

...

$$A_{m,1}x_1 + A_{m,2}x_2 + \cdots + A_{m,n}x_n = b_m$$

- If $m = n$ and \mathbf{A}^{-1} exists, then there is exactly one solution \mathbf{x}^* to the above equations
- If $m < n$, there are infinite number of solutions underdetermined
- If $m > n$, the solution may not exists! overdetermined

Norms

- Norm is used to measure the “size” of a vector or matrix
- The L_p norm of a vector is defined as

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{1/p}$$

for $p \in \mathbb{R}, p > 1$

- The L_2 norm is known as the **Euclidean norm**, and $\|x\|_2 = \sqrt{x^\top x}$
- The L_1 norm is simplified to $\|x\|_1 = \sum_i |x_i|$
- The L_∞ norm is simplified to $\|x\|_\infty = \max_i |x_i|$

Norms

- A norm is any function f that satisfies
 - $f(\mathbf{x}) \geq 0$
 - $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = 0$
 - $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (the triangle inequality)
 - $f(\alpha \mathbf{x}) = |\alpha|f(\mathbf{x}), \forall \alpha \in \mathbb{R}$
- Frobenius norm of a matrix

$$\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$$

Math basics

PROBABILITY THEORY

Random variable

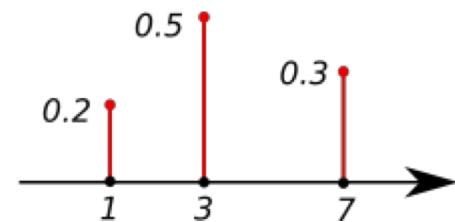
- A **random variable** is a variable that can take different values randomly
 - Denote the random variable by x and its two possible values by x_1 and x_2
 - For vectors, we write the random variable as \mathbf{x} and one of its values as x
 - Discrete versus continuous



Probability distribution

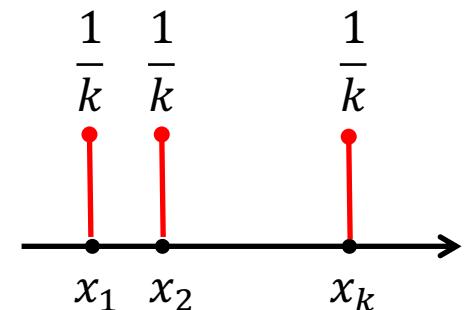
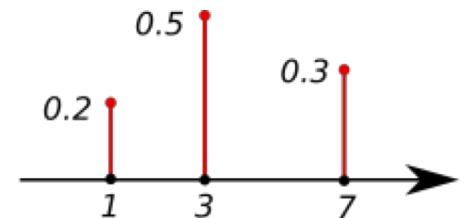
A **probability distribution** is a distribution of how likely a random variable or a set of random variables is to take on each of its possible states

- A probability distribution over **discrete** variables may be described using a **probability mass function (PMF)**
 - The prob that $x = x$ is denoted as $P(x)$ or $P(x = x)$
 - $x \sim P(x)$ specify which distribution x follows
- Joint probability
 - $P(x = x, y = y)$ or $P(x, y)$ denotes the prob that $x = x$ and $y = y$ simultaneously



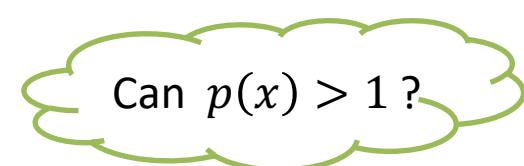
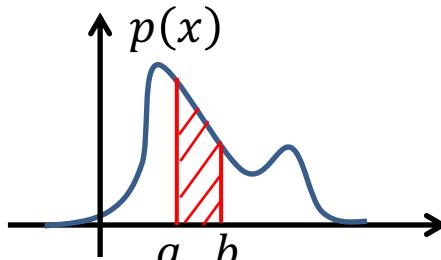
Probability mass function

- To be a PMF of a random variable x , a function P must satisfy:
 - The domain of P must be the set of all possible states of x
 - $\forall x \in X, 0 \leq P(x) \leq 1$
 - $\sum_{x \in X} P(x) = 1$
- Uniform distribution
 - Consider a single discrete random variable x with k different states
 - $P(x = x_i) = \frac{1}{k}, \forall i$



Probability density function

- A probability distribution over **continuous** variables may be described using a **probability density function** (PDF)
- To be a PDF, a function p must satisfy the following properties
 - The domain of p must be the set of all possible states of x
 - $\forall x \in \mathcal{X}, p(x) \geq 0$
 - $\int p(x)dx = 1$
- The prob that x lies in the interval $[a, b]$ is given by $\int_a^b p(x)dx$
- Note $p(x)$ does not give the prob of a specific state directly



Marginal probability

Suppose we know the prob distribution over a set of variables. The prob distribution over just a subset of them is known as the **marginal prob distribution**

- Let $P(x, y)$ denote the prob distribution of discrete random variables x and y , then

$$P(x = a) = \sum_b P(x = a, y = b)$$

- Let $p(x, y)$ denote the PDF of continuous random variables x and y , then

$$p(x) = \int P(x, y) dy$$

Bayes' rule



Thomas Bayes (1702~1761)

$$P(x|y) = \frac{P(x)P(y|x)}{P(y)}$$

Diagram illustrating the components of Bayes' rule:

- Prior probability of x (top left, blue text)
- Posterior probability of x (middle left, blue text)
- $P(x|y) = \frac{P(x)P(y|x)}{P(y)}$ (center equation)
- Probability of y given x ; likelihood (middle right, blue text)
- Prior probability of y (bottom right, blue text)

Conditional probability

The **conditional probability** is the probability of some event, given that some other event has happened

- The conditional prob that $y = b$ given $x = a$ is denoted by $P(y = b|x = a)$, which can be calculated as

$$P(y = b|x = a) = P(y = b, x = a)/P(x = a)$$

- The chain rule

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)}|x^{(1)}, \dots, x^{(i-1)})$$

- **Exercise:** With this definition, is the following correct?

$$P(a, b, c) = P(a|b, c)P(b|c)P(c)$$

Independence and conditional independence

Two random variables x and y are **independent** if their joint prob distribution can be expressed as a product of two factors, one involving x and one involving y :

$$\forall x \in X, y \in Y,$$

$$p(x = a, y = b) = p(x = a)p(y = b)$$

Denoted by $x \perp y$

Two random variables x and y are **conditionally independent** given a random variable z if the conditional prob distribution over x and y factorizes in this way:

$$\forall a \in X, b \in Y, c \in Z,$$

$$p(x = a, y = b | z = c) = p(x = a | z = c)p(y = b | z = c)$$

Denoted by $x \perp y | z$

Expectation

The **expectation**, or **expected value**, of some function $f(x)$ w.r.t. a prob distribution $P(x)$ is the average value that f takes on when x is drawn from P

- For discrete variables

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x)$$

- For continuous variables

$$\mathbb{E}_{x \sim P}[f(x)] = \int P(x)f(x)dx$$

- If the identity of the distribution is clear, we may write $\mathbb{E}_x[f(x)]$
- Expectation is **linear**: if α and β do not depend on x , then

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)]$$

Variance and covariance

- The **variance** of a function $f(x)$

$$\text{Var}(f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

- It measures how much the value of the function f varies when x is sampled from its distribution

- The **covariance** of two functions $f(x)$ and $g(x)$

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$

- It measures how much two values are linearly related to each other, as well as the scale of these variables
 - High absolute values of the covariance mean that the values change very much and are both far from their respective means
- The **covariance matrix** of a random vector $\mathbf{x} \in \mathbb{R}^n$ is an $n \times n$ matrix

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(\mathbf{x}_i, \mathbf{x}_j)$$

Common prob distributions

- Bernoulli distribution: over a single binary random variable

$$P(x = 1) = \phi, P(x = 0) = 1 - \phi$$

$$P(x = x) = \phi^x(1 - \phi)^{1-x}$$

$$\mathbb{E}_x[x] = \phi, \text{Var}(x) = \phi(1 - \phi)$$

- Multinoulli or categorical distribution: over a single discrete variable with k diff states where k is finite

$$P(x = i | \mathbf{p}) = p_i$$

– where $\mathbf{p} \in [0,1]^k$ and $\sum_{i=1}^k p_i = 1$

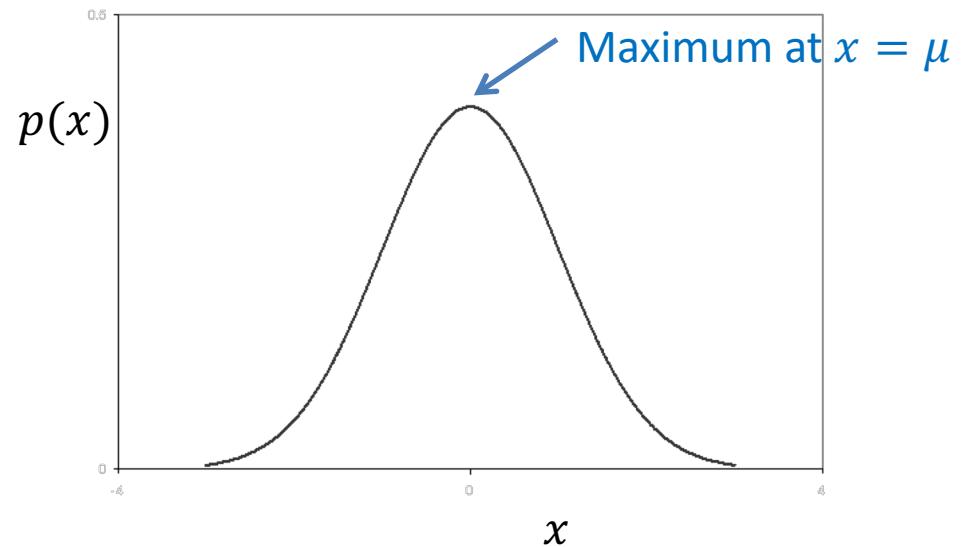
Common prob distributions

- Gaussian distribution or normal distribution: over a continuous variable

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- Mean: $\mathbb{E}[x] = \mu$
- Variance: $\text{Var}[x] = \sigma^2$
- Standard deviation: σ

The *central limit theorem* shows that the sum of many independent variables is approximately normally distributed



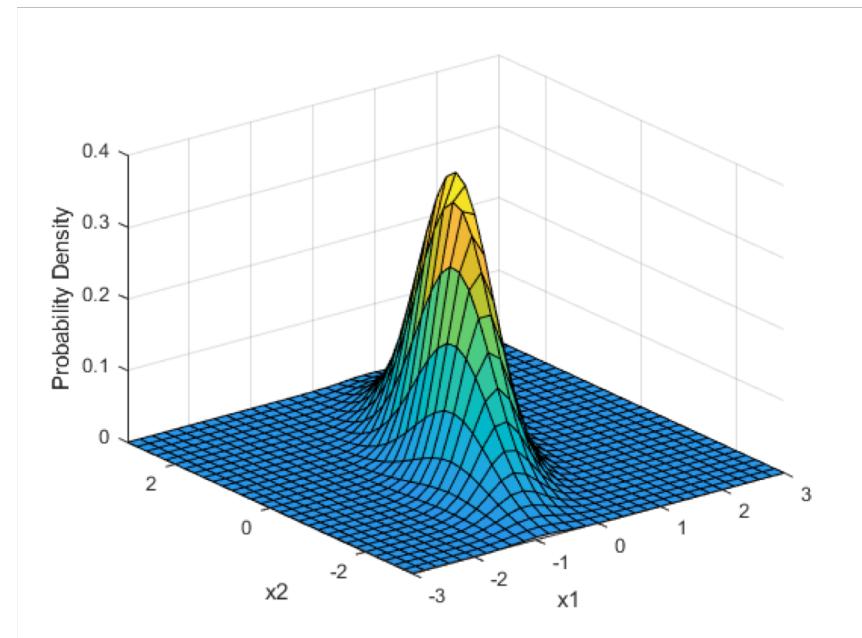
Common prob distributions

- Multivariate normal distribution: over a continuous vector

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{(2\pi)^2 \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

- Mean: $\mathbb{E}[x] = \mu$
- Covariance: $\text{Cov}[x] = \Sigma$

Isotropic Gaussian distribution: $\Sigma = \alpha I$
where α is a scalar



Math basics

OPTIMIZATION

Gradient-based optimization

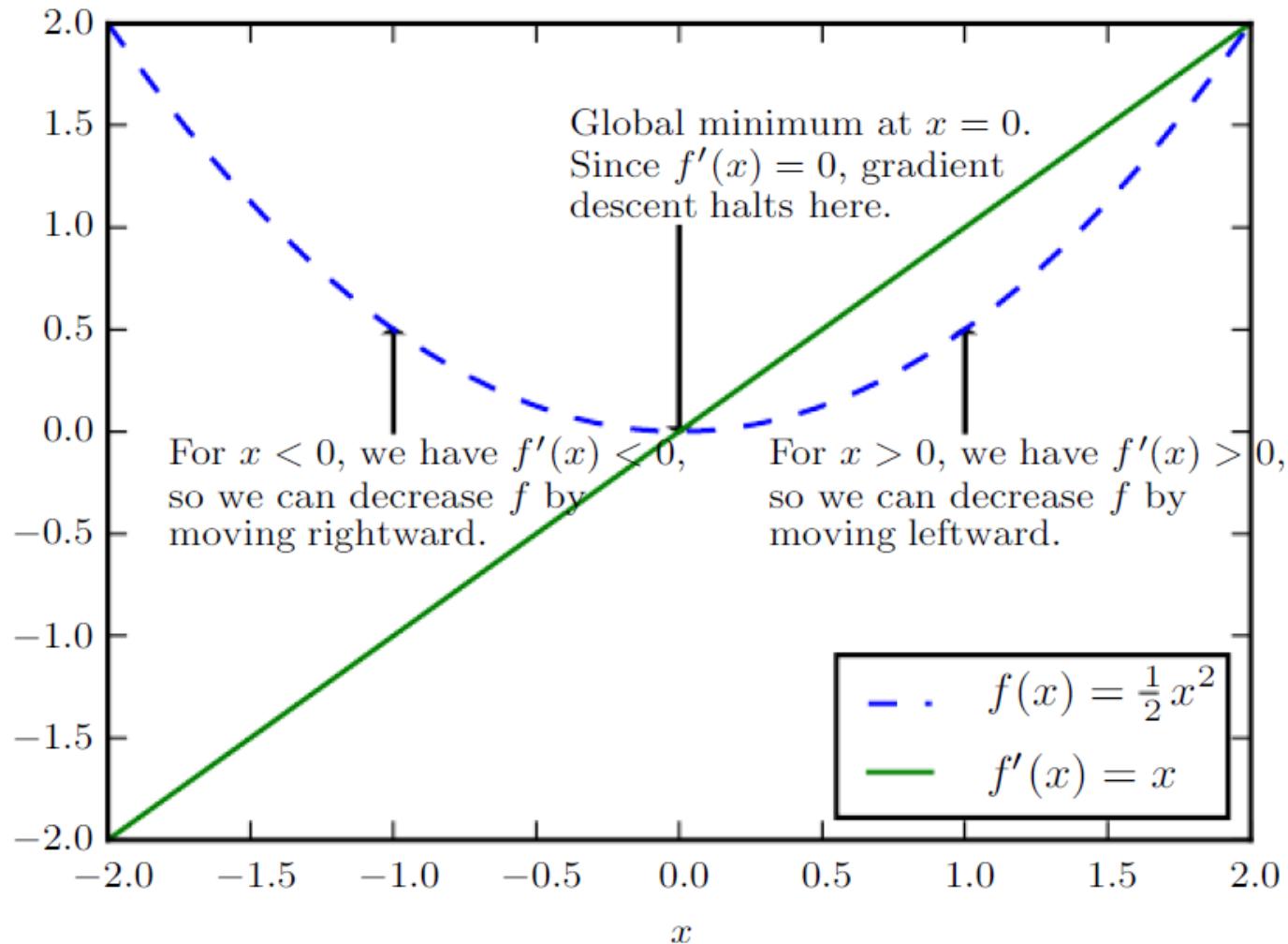
- The function we want to minimize or maximize is called **objective function**
- When we are minimizing it, we may also call it the **cost function**, **loss function**, or **error function**
- The **derivative** of a function $y = f(x)$, denoted by $f'(x)$ or $\frac{dy}{dx}$, gives the slope, or **gradient**, of f at the point x
- Gradient descent

$$x' = x - \eta f'(x)$$

where $\eta > 0$ is the **learning rate**

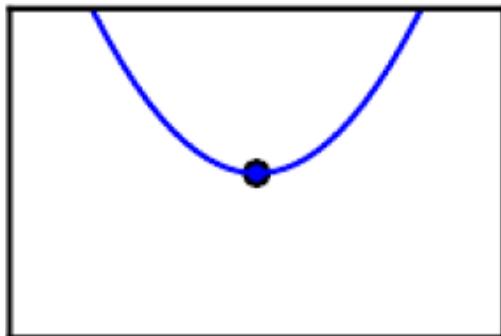
Gradient descent

$$x' = x - \eta f'(x)$$

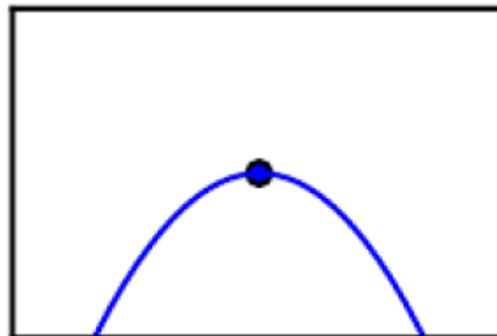


Critical points

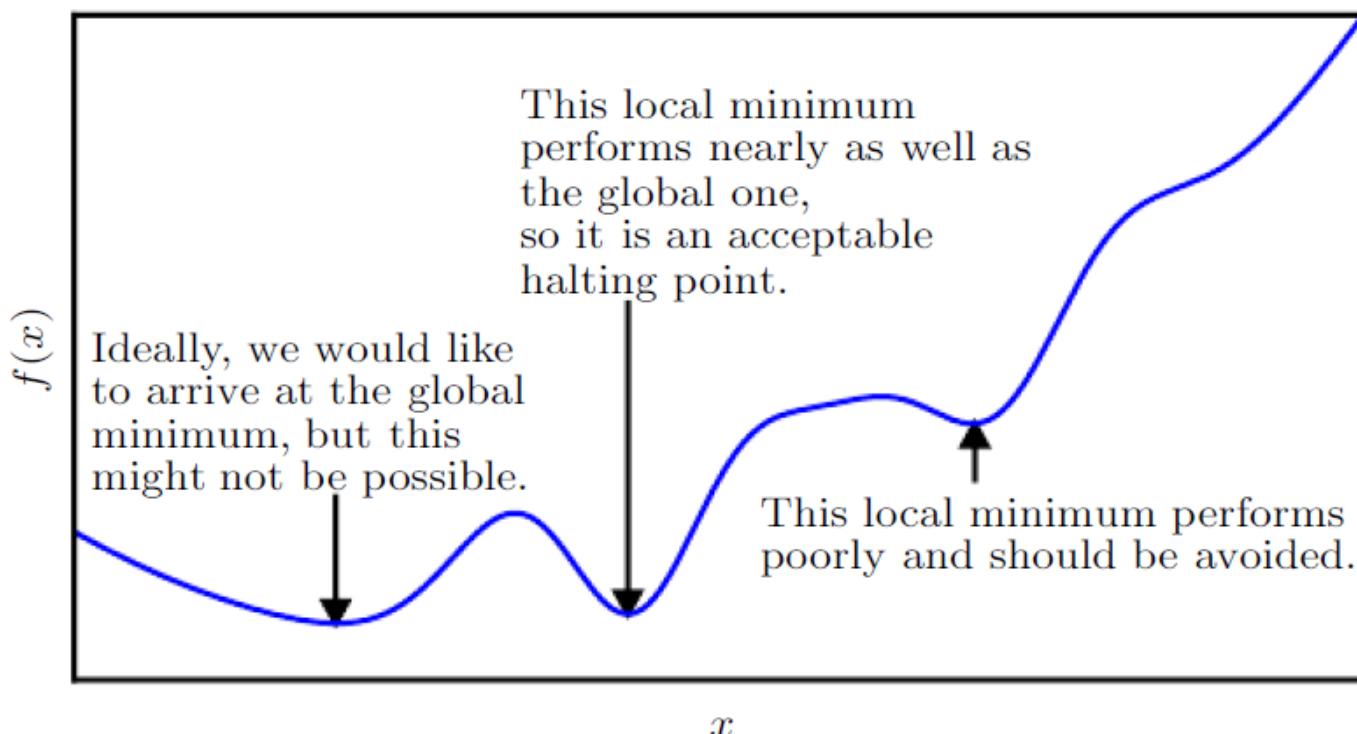
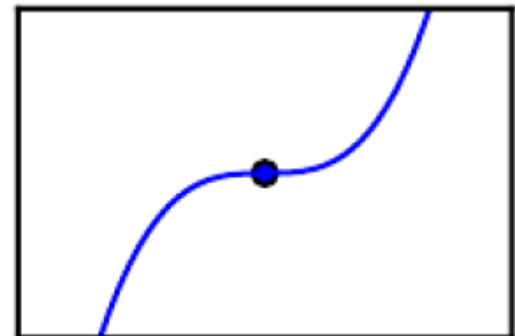
Minimum



Maximum



Saddle point



Gradient decent for multivariate functions

- For a function of a single variable $y = f(x)$, the gradient decent method is

$$x' = x - \eta f'(x)$$

- For a function $y = f(\mathbf{x})$, the **partial derivative** is denoted by $\partial f / \partial x_i$
- The gradient decent method becomes

$$\mathbf{x}' = \mathbf{x} - \eta \nabla_{\mathbf{x}} f(\mathbf{x})$$

where $\eta > 0$ and $\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{pmatrix}$

Summary so far

- Linear algebra
 - Math objects: Scalars, vectors, matrices, tensors
 - Simple operations: matrix transpose, inverse, product
 - Norms: L_p norm
- Probability theory
 - Random variables: discrete, continuous
 - Prob distribution: PMF and PDF
 - Marginal probability
- Conditional probability
- Independence and conditional independence
- Expectation, variance and covariance
- Common prob distributions
- Bayes' rule
- Optimization
 - Gradient descent
 - Critical points

Outline

- Overview of deep learning
- Math basics
- Machine learning basics

Learning algorithms

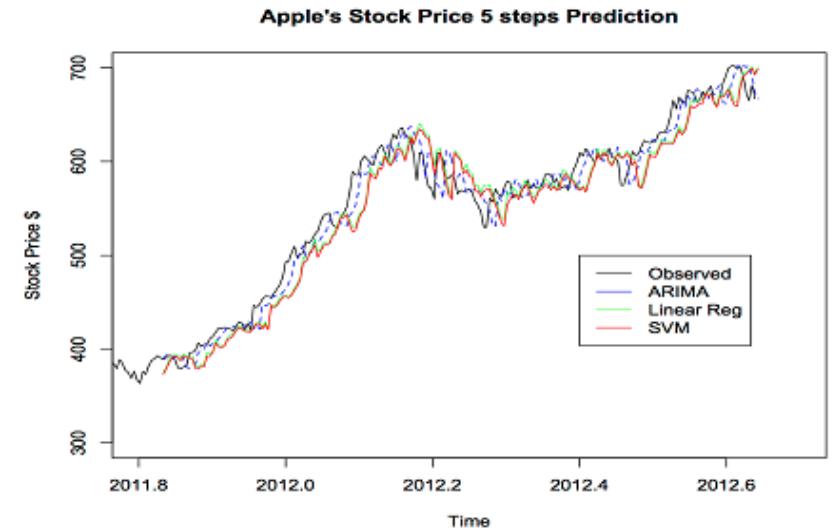
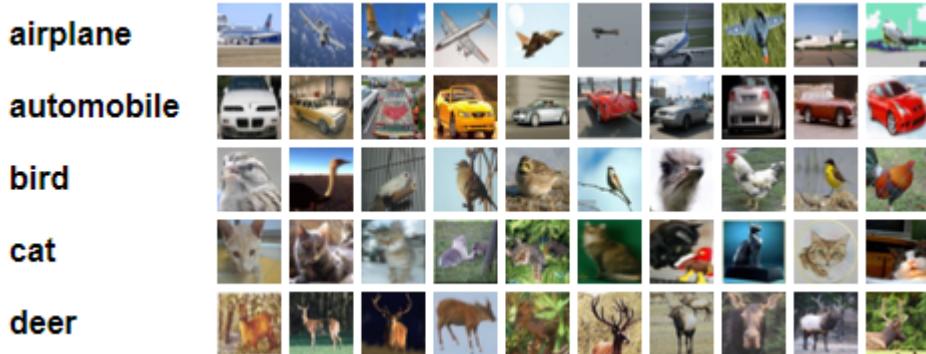
“A computer program is said to **learn** from experience E w.r.t. some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .” ---Tom Mitchell, 1997

- Machine learning (ML) tasks are usually described in terms of how the ML system should process an **example**
- An example is a collection of **features** that have been quantitatively measured from some object or event
 - Features of a bucket: color, diameter, height, material, etc
 - Features of an animal: size, shape, number of legs, , etc



The tasks T

- Classification
 - Suppose there are k categories. Find a function $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$



- Regression
 - Find a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, and m is often 1

Regression results might be converted to classification results, if the regression curves have special forms

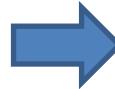
The tasks T

- Synthesis and sampling dataset

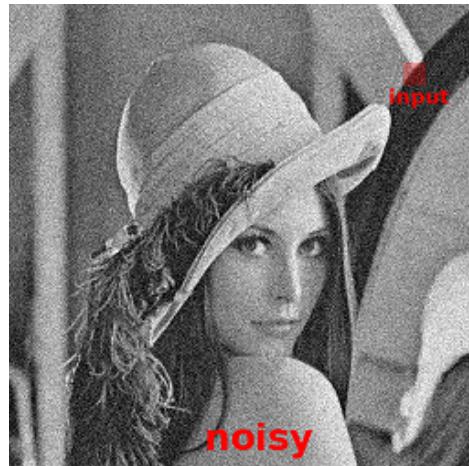
5 0 4 1 9 2 1 3 1 4
3 5 3 6 1 7 2 8 6 9
4 0 9 1 1 2 4 3 2 7
3 8 6 9 0 5 6 0 7 6
1 8 7 9 3 9 8 5 9 3
3 0 7 4 9 8 0 9 4 1
4 4 6 0 4 5 6 1 0 0

Synthesized using GAN

7	3	9	3	9
1	1	0	6	0
0	1	9	1	2
6	3	2	0	8

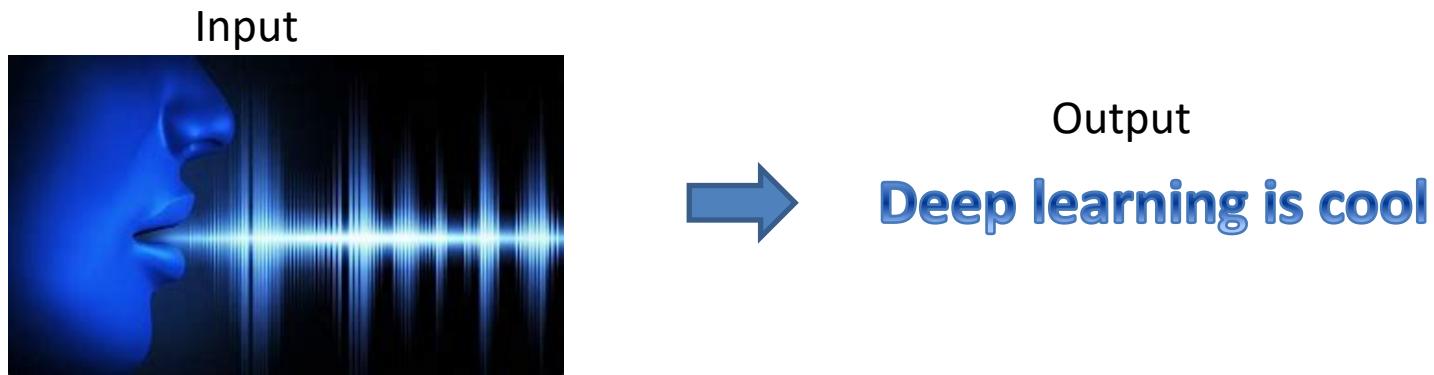


- Denoising



The tasks T

- Transcription



- Machine translation



The tasks, T

- Structured output
- Anomaly detection
- Synthesis and sampling
- Imputation of missing values
- Density estimation
- Etc.

The performance measure, P

- A performance measure is required to quantitatively evaluate the performance of a ML system
- Usually this measure P is **specific to the task T** being carried out by the system
 - Classification and transcription: accuracy or error rate
 - Regression and denoising: distance between the ground-truth and prediction
 - Synthesis, machine translation: difficult and sometimes need human evaluation
- What we are more interested in is the performance measure on a **test set** of data that is **separated** from the data used for training the system

The experience, E

- ML algorithms including deep learning algorithms can be broadly categorized as **unsupervised** or **supervised** by what kind of experience they are allowed to have during the learning process
- The algorithms experience a **dataset**, which is a collection of many **examples** or **data points** denoted by x
 - We can view examples as samples of a random variable \mathbf{x}
- Unsupervised learning algorithms
 - experience a dataset containing many features, then learn useful properties or structures of this dataset
- Supervised learning algorithms
 - experience a dataset containing many features, but each example is also associated with a **label** or **target** denoted by y

learn $p(\mathbf{x})$

learn $p(y|\mathbf{x})$

Example: linear regression

- Task T : to predict y from x by outputting $\hat{y} = \mathbf{w}^\top \mathbf{x}$ x_i : feature
 w_i : weight
- Performance P : mean squared error of the model on the test with m test samples $\{(\mathbf{x}_i, y_i)\}^{\text{test}}$

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{y}_i - y)^{\text{test}}$$

- Experience E : minimize the MSE on the training set of q samples $\{(\mathbf{x}_i, y_i)\}^{\text{train}}$

$$\text{MSE}_{\text{train}} = \frac{1}{q} \sum_i (\hat{y}_i - y)^{\text{train}}$$

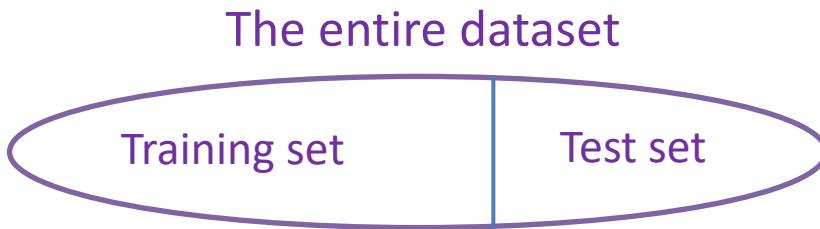
- Denote $\{(\mathbf{x}_i, y_i)\}^{\text{train}}$ collectively by $(\mathbf{X}^{\text{train}}, \mathbf{y}^{\text{train}})$, then

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{w}} \frac{1}{q} \left\| \hat{\mathbf{y}}^{\text{train}} - \mathbf{y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^{\text{train}}{}^\top \mathbf{X}^{\text{train}})^{-1} \mathbf{X}^{\text{train}}{}^\top \mathbf{y}^{\text{train}}$$

Capacity, overfitting and underfitting

- A ML algorithm must perform well on new, previously unseen inputs—not just on which it was trained
 - This ability is called generalization



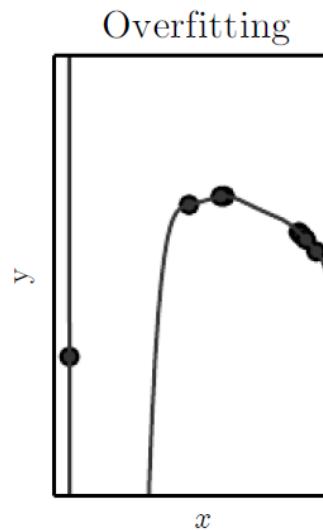
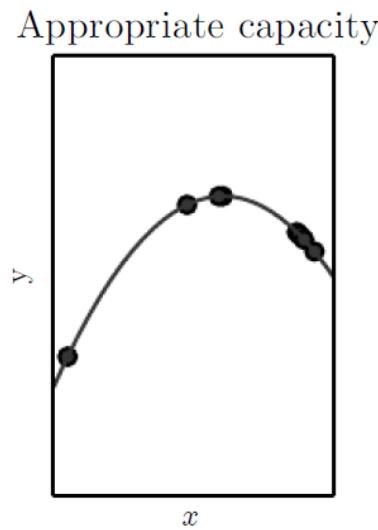
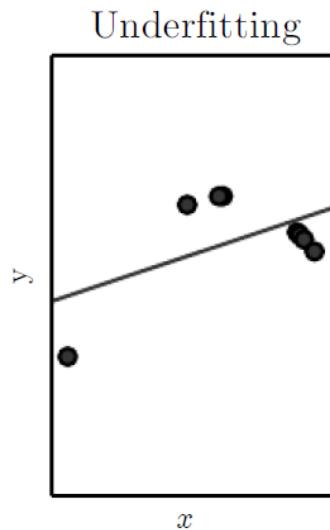
- Smaller training error → higher model capacity
 - If the training error is too large, the model is underfitting the training set
- Smaller test error or generalization error → higher generalization ability
 - If the training error is very small but the test error is very large, the model is overfitting the training set

Example: polynomial regression

- Consider a regression problem in which the input x and output y are both scalars. Find a function $f: \mathbb{R} \rightarrow \mathbb{R}$ to fit the data
 - $f(x) = b + wx$
 - $f(x) = b + w_1x + w_2x^2$
 - $f(x) = b + \sum_{i=1}^9 w_i x^i$

MSE training:

$$\min_w \frac{1}{N} \sum_{n=1}^N \|f(x^{(n)}) - y^{(n)}\|_2^2$$



General principles

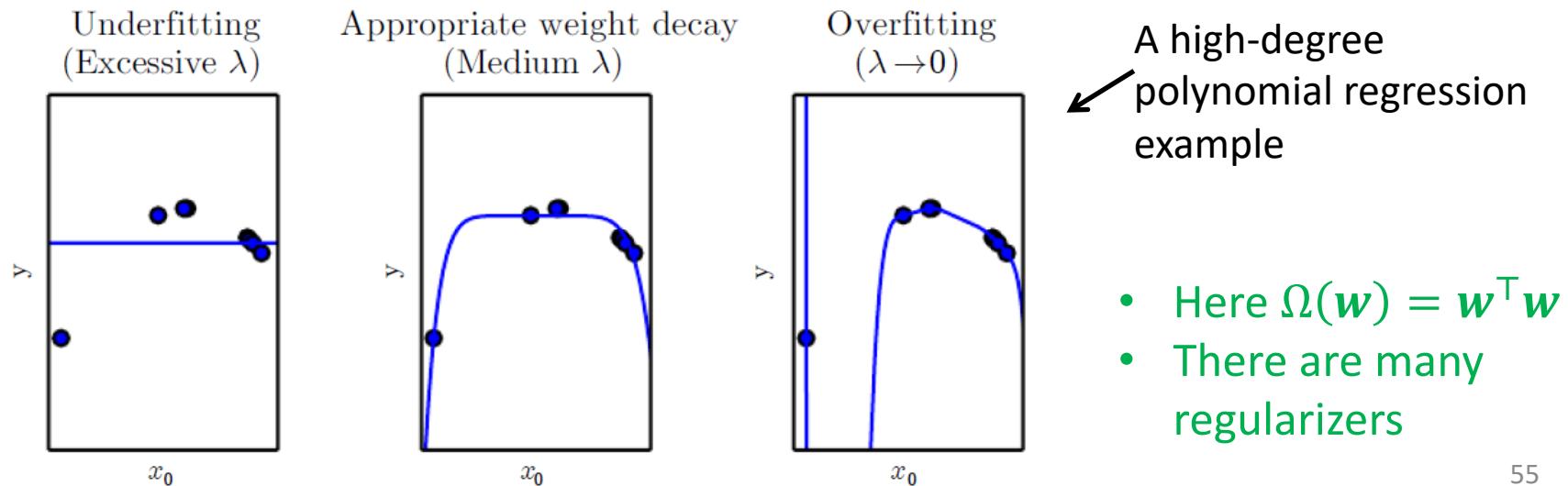
- Increase the model capacity
 - Make the training error small
- Increase the generalization ability
 - Make the gap between training error and test error small

Regularization

- To carry out a specific task, we often build a set of **preferences** into the learning algorithm, which is embodied by a **regularizer** Ω
- E.g., for polynomial regression, the total cost function becomes

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

where $\lambda > 0$ is a constant.



Regularization

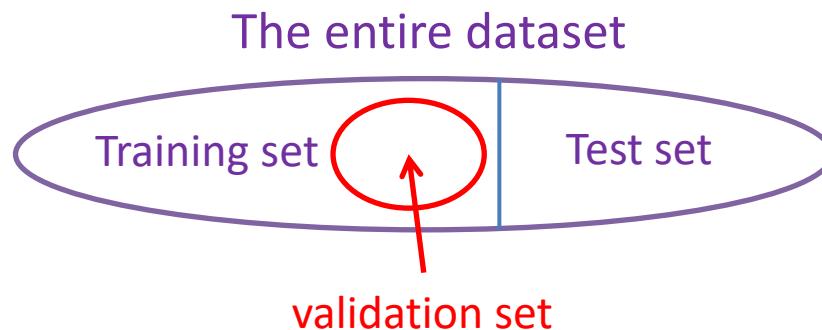
Regularization is any modification we make to a learning algorithm that is intended to reduce its **generalization error** but not its training error

Hyperparameters

- Many machine learning algorithms have two sets of parameters:
 - **Hyperparameters**: control the algorithm's behavior and are not adapted by the algorithm itself. They often determine the **capacity** of the model
 - **Learnable parameters** (“learnable” is often omitted): can be learned from data
- The polynomial regression algorithm $J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$
 - Hyperparameters: λ
 - Learnable parameters: \mathbf{w}
- A neural network
 - Hyperparameters: the number of layers, the number of neurons per layer, etc.
 - Learnable parameters: weights and bias in each layer

Validation sets

- How to choose the hyperparameters considering that we cannot see the test set?
 - Set them such that the training error is as small as possible?
- We need another set on which the model is not trained on
 - Make the error on this set as small as possible
 - This is called the **validation set**
- How do we obtain a validation set?



Maximum likelihood estimation

Problem definition

- Given a set of m examples $\mathbb{X} = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(\mathbf{x})$
 - Find a prob distribution $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ to approximate $p_{\text{data}}(\mathbf{x})$
 - The task is to find optimal $\boldsymbol{\theta}$
-
- The maximum likelihood estimator for $\boldsymbol{\theta}$ is defined as
$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$
 - We usually use

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})\end{aligned}$$

- where \hat{p}_{data} is the empirical distribution

Conditional log-likelihood

- Estimate a conditional probability $P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$ in order to predict \mathbf{y} given \mathbf{x}
 - E.g. For classification, \mathbf{y} is a (discrete) random variable representing label of an input \mathbf{x}
- If X represents all inputs and Y all observed targets, then the **conditional maximum likelihood estimator** is
$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P_{\text{model}}(Y|X; \boldsymbol{\theta})$$
- If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Stochastic gradient decent



- Minimizing the cost function over the entire training set is computationally expensive
- We often decompose the training set into smaller subsets or **minibatches** and optimize the cost function defined over individual minibatches $(\mathbf{X}^{(i)}, \mathbf{y}^{(i)})$ and take the average

$$J(\boldsymbol{\theta}) = \frac{1}{m'} \sum_{i=1}^{m'} L(\mathbf{X}^{(i)}, \mathbf{y}^{(i)}, \boldsymbol{\theta})$$

$$\mathbf{g} = \frac{1}{m'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{m'} L(\mathbf{X}^{(i)}, \mathbf{y}^{(i)}, \boldsymbol{\theta})$$

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \mathbf{g}$$

- A total of m' minibatches
- The batchsize ranges from 1 to a few hundreds

Summary

- Overview of deep learning
- Math basics
 - Linear algebra
 - Probability theory
 - Optimization
- Machine learning basics
 - Learning algorithms: Task, performance measure, experience
 - Capacity, Overfitting and Underfitting
 - Hyperparameters and Validation Sets
 - Maximum Likelihood Estimation
 - Stochastic gradient decent