

Deep Learning in Computer Vision

王鹏

peng.wang@nwpu.edu.cn

Pre-requisite

- Elementary programming ability
 - Most of assignments may be written in python
 - Suggest Pytorch as your deep learning platform
- Matrix/Vector Calculus, Linear Algebra, Probability and Statistics

Books and References

- Dive into Deep Learning, Aston Zhang and Zachary C. Lipton and Mu Li and Alexander J. Smola, <http://www.d2l.ai>, 2020.
- Deep Learning, Ian Goodfellow and Yoshua Bengio and Aaron Courville, MIT Press, 2016
- Pattern Recognition and Machine Learning, Christopher Bishop, Springer, 2011
- The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition, Springer, 2nd edition, 2016
- Convex Optimization, Stephen Boyd and Lieven Vandenberghe, Cambridge University Press, 2004

Grading Policy

- Three assignments: **15% x 3 = 45%**
- Final exam: **45%**
- In-class Participation: **10%**

What we will learn

- Basic concepts about artificial intelligence, machine learning, computer vision, numeric optimization
- Deep learning: Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), Attention
- The application of DL to CV tasks, such as image classification, object detection, image segmentation, image caption and visual question answering

What we will miss

- We focus on data-driven AI methods, so will ignore symbolic knowledge-driven methods
- We focus on discriminative ML algorithms, so will ignore generative ML algorithms such as graphic models
- We focus on supervised DL methods, so will ignore unsupervised and semi-supervised DL methods such as autoencoders
- We focus on mid-to-high level CV tasks, so will ignore low level tasks such as image denoise

Schedule and Syllabus

L1	Introduction	L11	DNNs for Pixel Prediction Tasks
L2	Linear Algebra Basics + Matrix Calculus	L12	DNNs for Object Detection (Two-stage, One-stage, Anchor-free)
L3	Optimization Basics	L13	DNNs for Object Tracking
L4	Machine Learning Basics	L14	DNNs for Vision-Language Tasks (CNN-LSTM, BERT, MCB)
L5	Neural Networks (Multi-layer Perceptrons, Activation Functions, Back-propagation)	L15	DNNs for Human-related Vision Tasks (Face Recognition, Pose Estimation and Person Re-identification)
L6	Training Neural Networks (SGD, Tricks for minimizing training errors and generalization gaps)	L16	DNNs for 3D Vision and DNNs for Action&Video
L7	Convolutional Neural Networks and Image Classification	L17	Graph Neural Networks for Visual Reasoning (GCN, GGNN, GAN)
L8	Recurrent Neural Networks	L18	AutoML for Vision Tasks (Classification, Segmentation, Detection)
L9	Attention Mechanisms (soft, hard, spatial, channel and self attentions)	L19	Visualization and Interpretability
L10	Generative Adversarial Networks	L20	Future Directions and Discussion

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms



Image Deblur

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms



Image Denoise

Computer Vision

Tasks



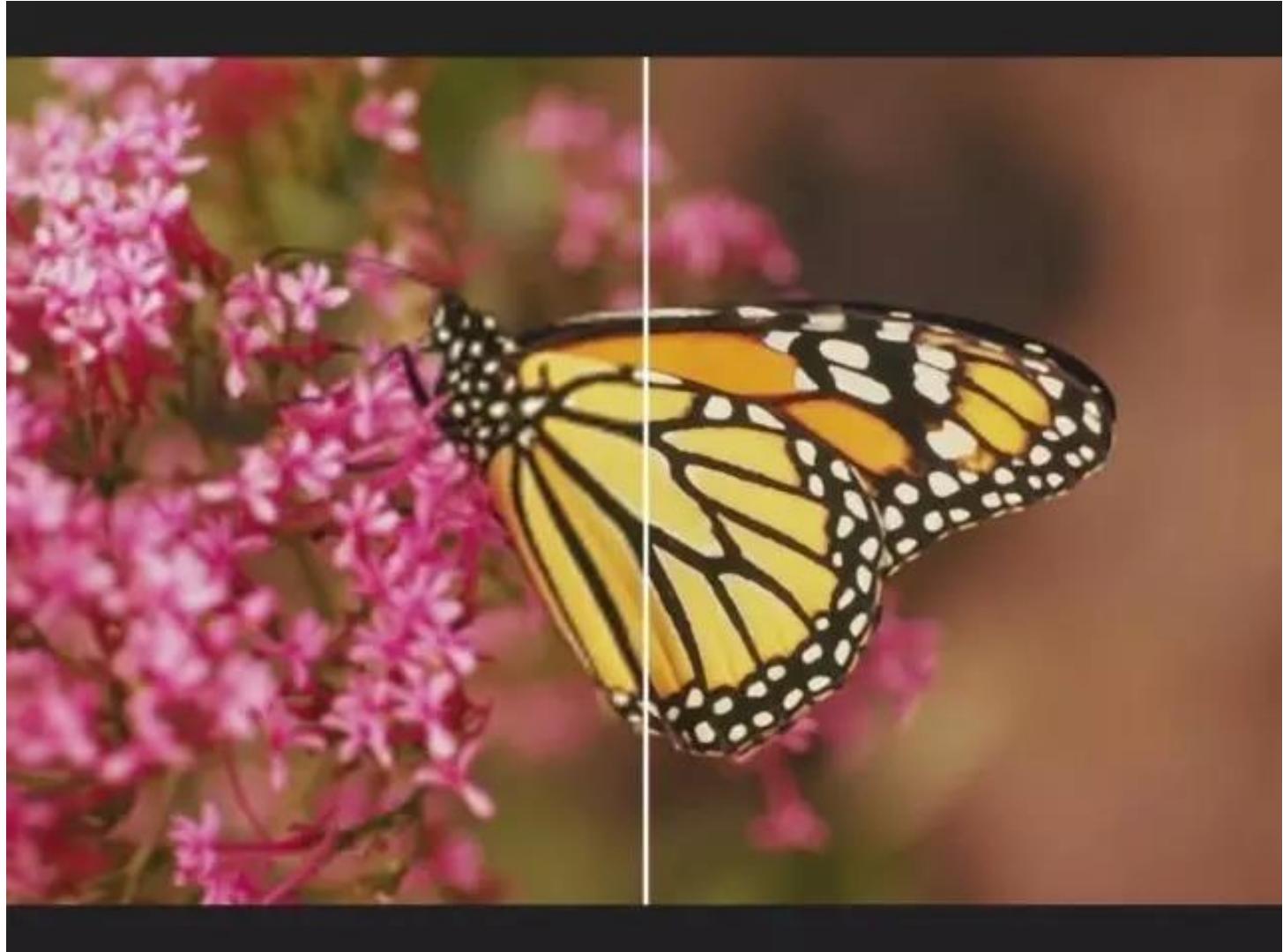
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Super-resolution

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms

Image Classification

Easiest classes

red fox (100) hen-of-the-woods (100) ibex (100) goldfinch (100) flat-coated retriever (100)



tiger (100)



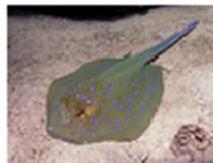
hamster (100)



porcupine (100)



stingray (100)



Blenheim spaniel (100)



Hardest classes

muzzle (71)



hatchet (68)



water bottle (68)



velvet (68)



loupe (66)



hook (66)



spotlight (66)



ladle (65)



restaurant (64)



letter opener (59)



Image Classification

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms



Image Segmentation

Computer Vision

Tasks



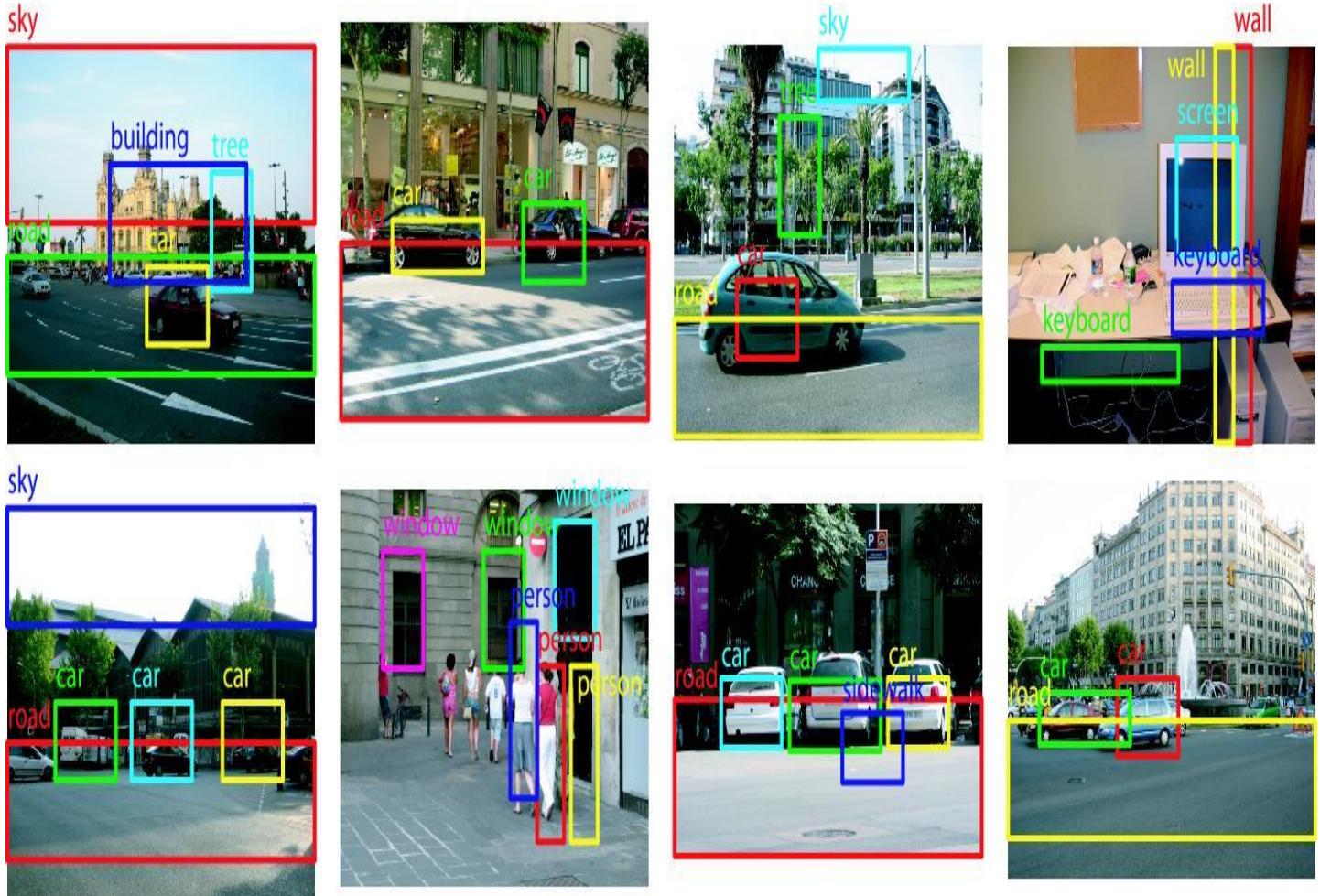
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Object Detection

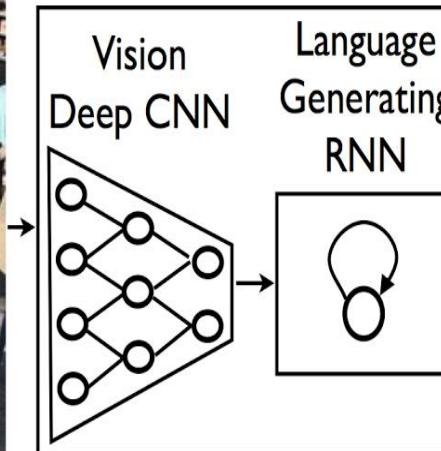
Computer Vision

Tasks

Datasets

Feature
Engineering

Machine
Learning
Algorithms



A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

Image Captioning

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms



What is the mustache
made of?

AI System

bananas

Visual Question Answering

Computer Vision

Tasks



Datasets



Feature
Engineering



Machine
Learning
Algorithms

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

MNIST (1998)

Computer Vision

Tasks



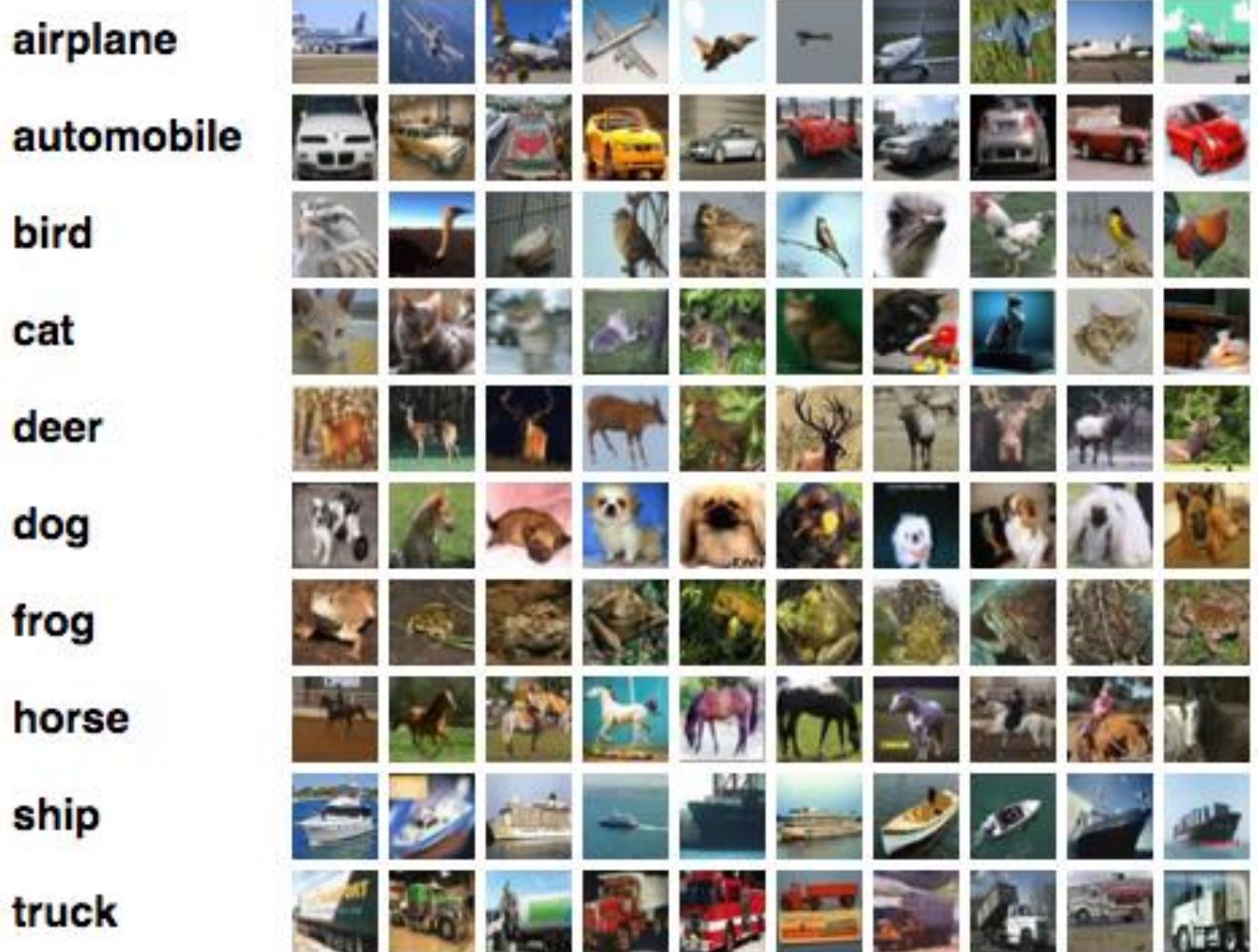
Datasets



Feature
Engineering



Machine
Learning
Algorithms



CIFAR10 & CIFAR100

Computer Vision

Tasks



Datasets



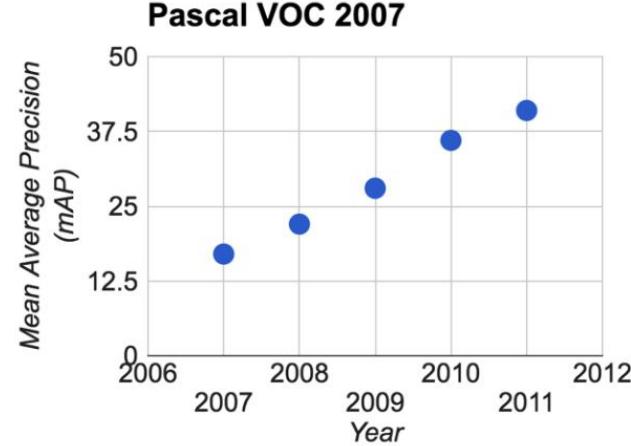
Feature
Engineering



Machine
Learning
Algorithms

PASCAL Visual Object Challenge (20 object categories)

[Everingham et al. 2006-2012]



Computer Vision

Tasks



Datasets



Feature
Engineering



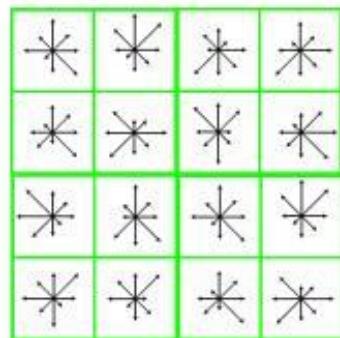
Machine
Learning
Algorithms

The SIFT Object Recognition Algorithm

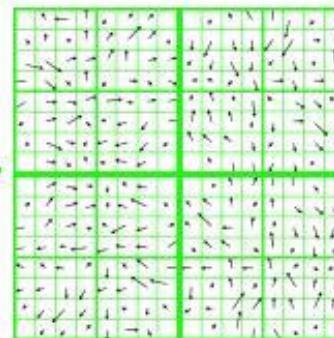
Incrementally Gaussian Blur
The Original Image to
Create a Scale
Space



Find the Difference Between
Adjacent Gaussian Images
in Scale Space

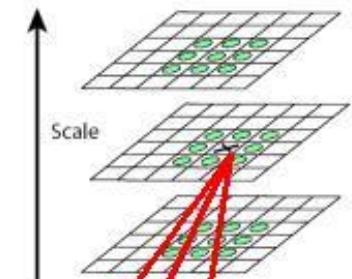


Sixteen Histograms are
Created Using The Gradients.
Using 8 Orientations, This
Makes 128-D Feature Vectors.



The Gradient of Pixels Around
Each Keypoint is Determined
At the Gaussian Scale at Which
It Was Found

Keypoints are Pixels
in Difference Images
That are Larger Than
or Smaller Than all 26
Neighbors



Hundreds of
Keypoints are Found

SIFT Feature (1999) <http://weitz.de/sift/index.html?size=large>

Computer Vision

Tasks



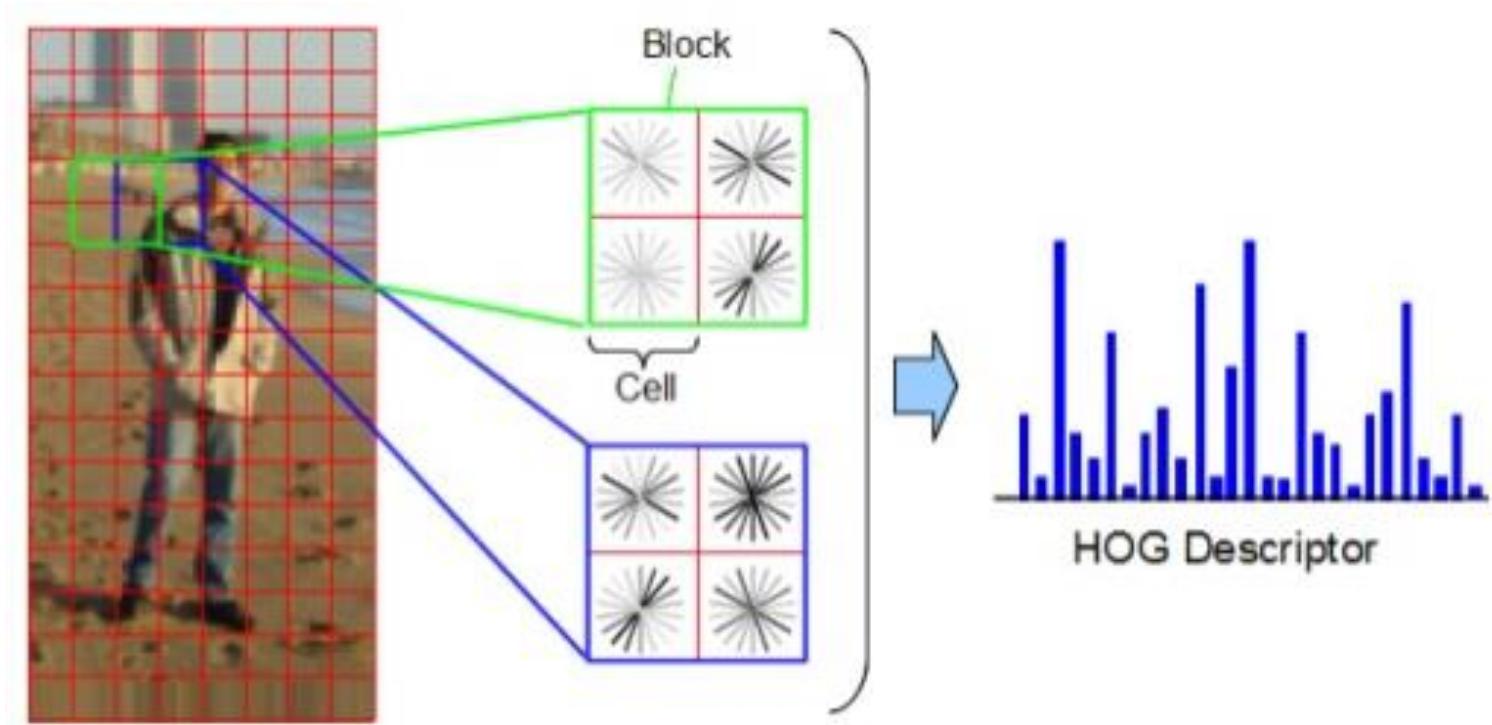
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Histogram of oriented gradients (2005)

Computer Vision

Tasks



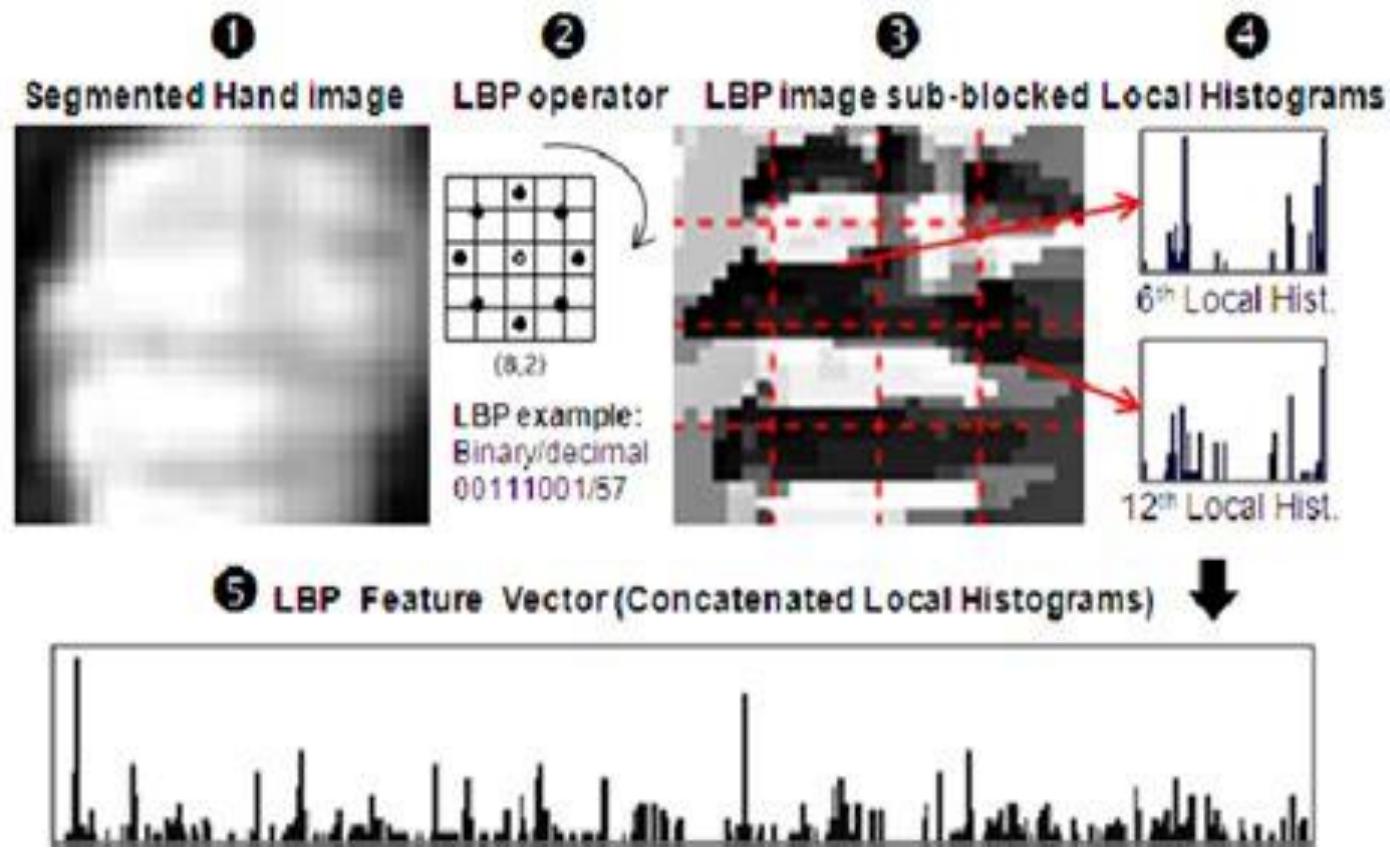
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Local Binary Patterns (1994)

Computer Vision

Tasks



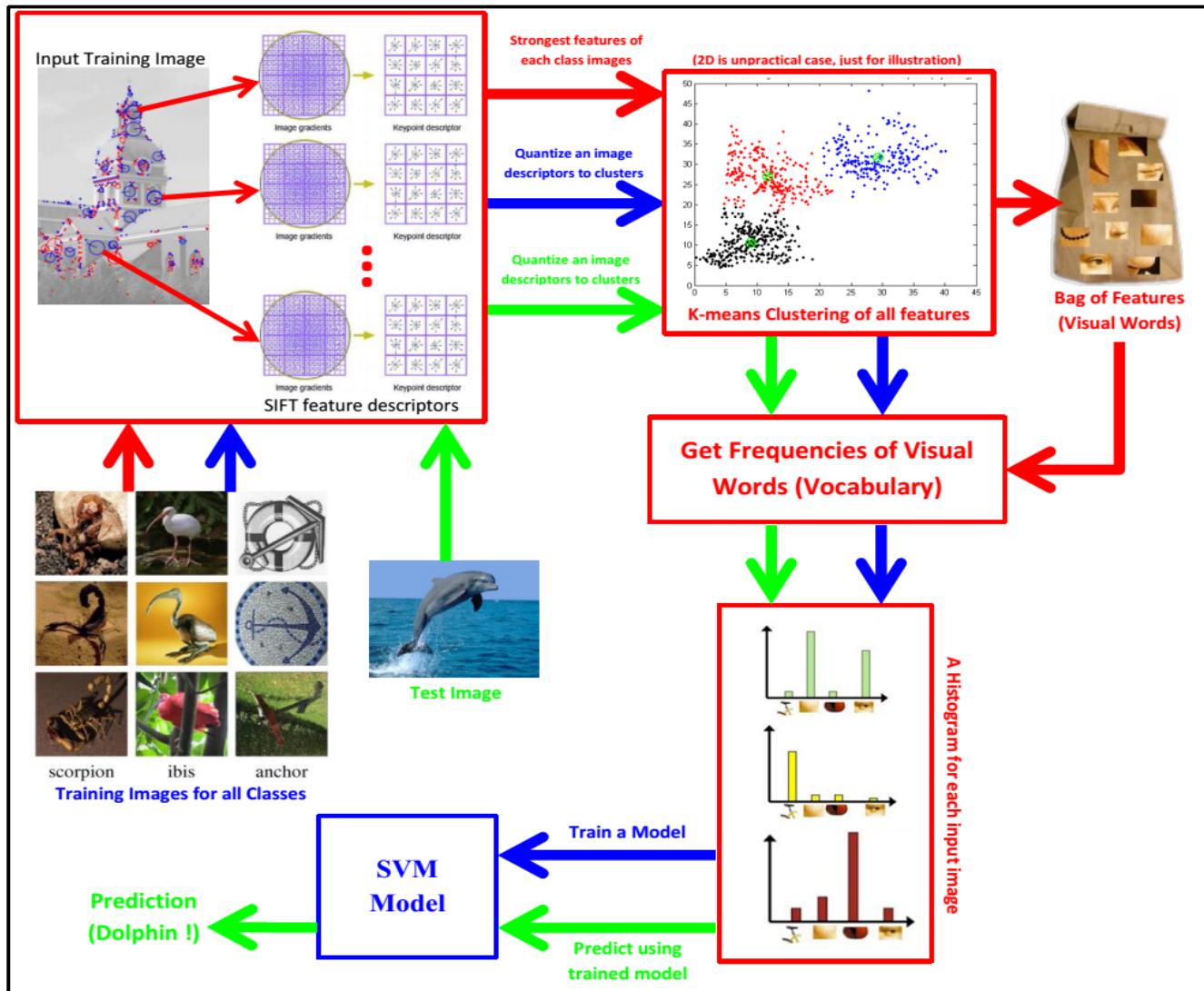
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Bag-of-visual-words for image classification

Computer Vision

Tasks



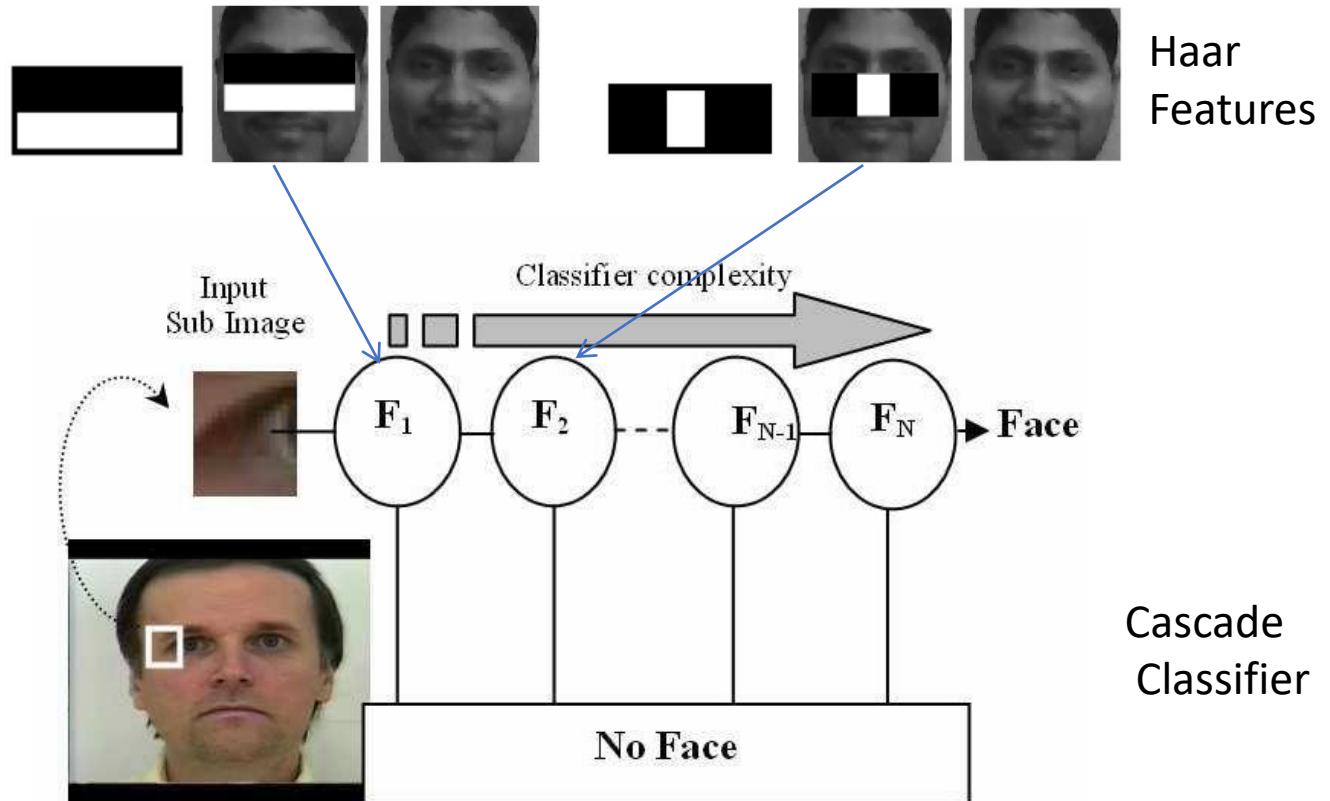
Datasets



Feature
Engineering



Machine
Learning
Algorithms



Viola-Jones Real-time face detector (2001)

Computer Vision

Tasks



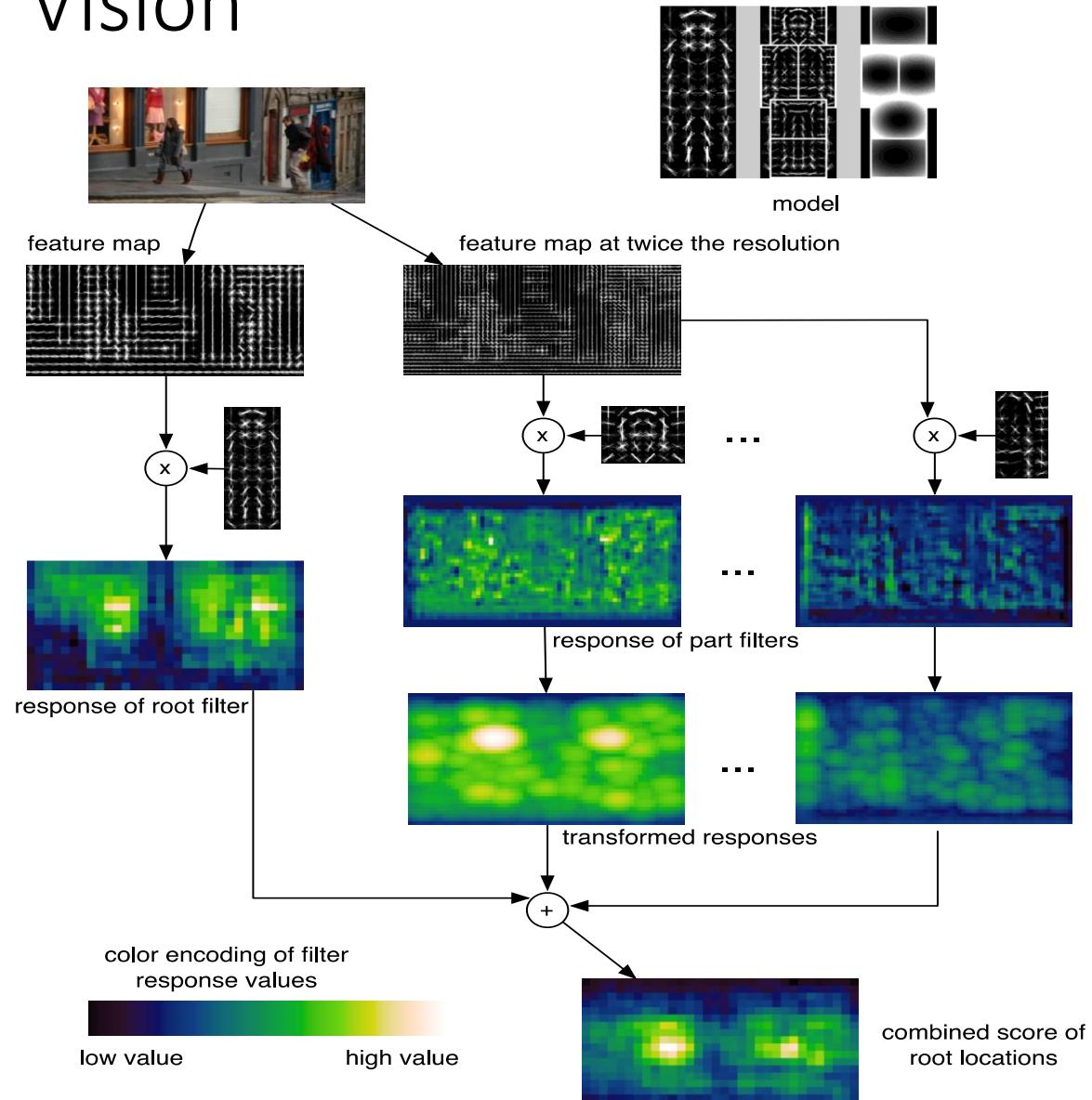
Datasets



Feature
Engineering

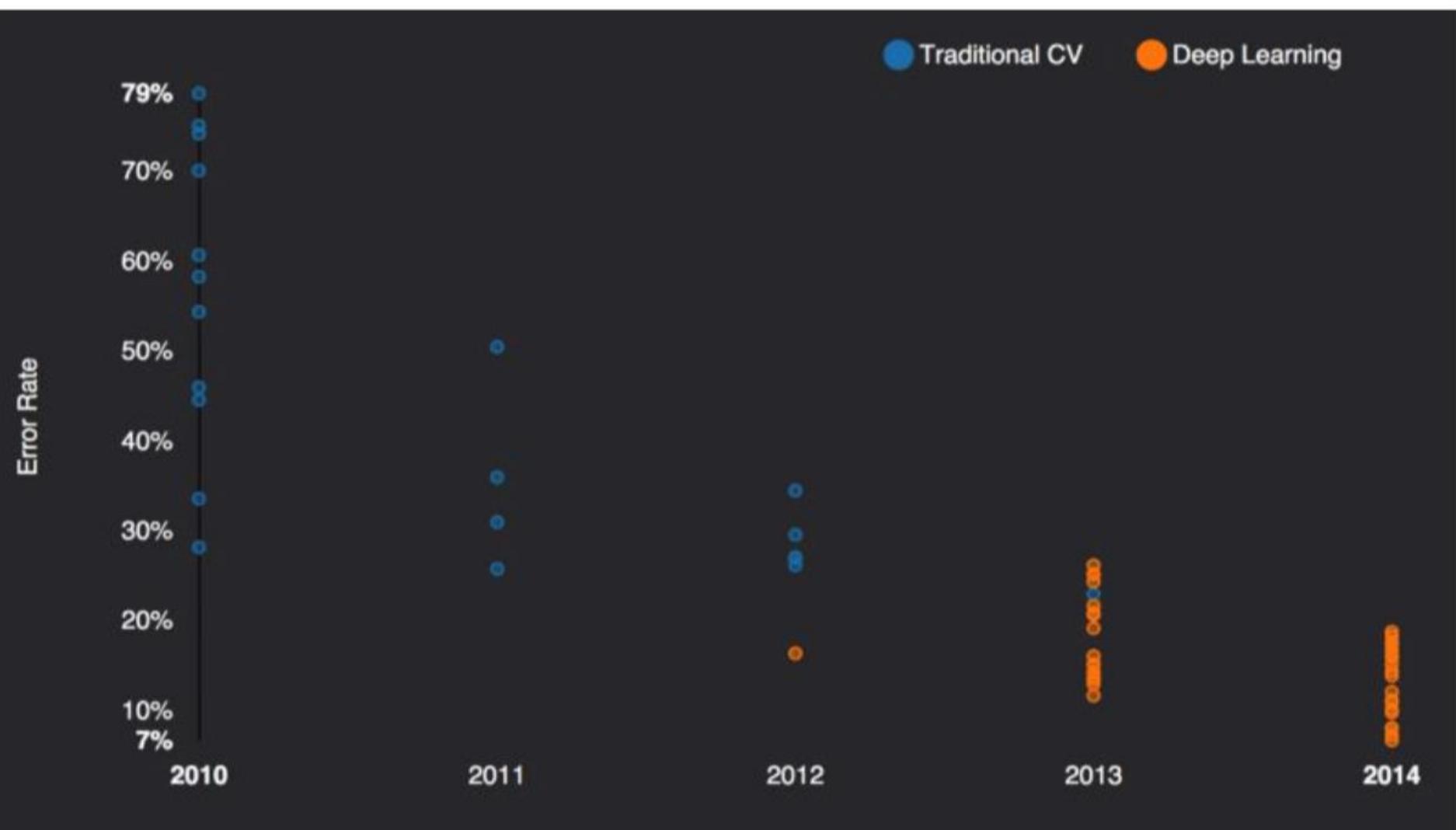


Machine
Learning
Algorithms



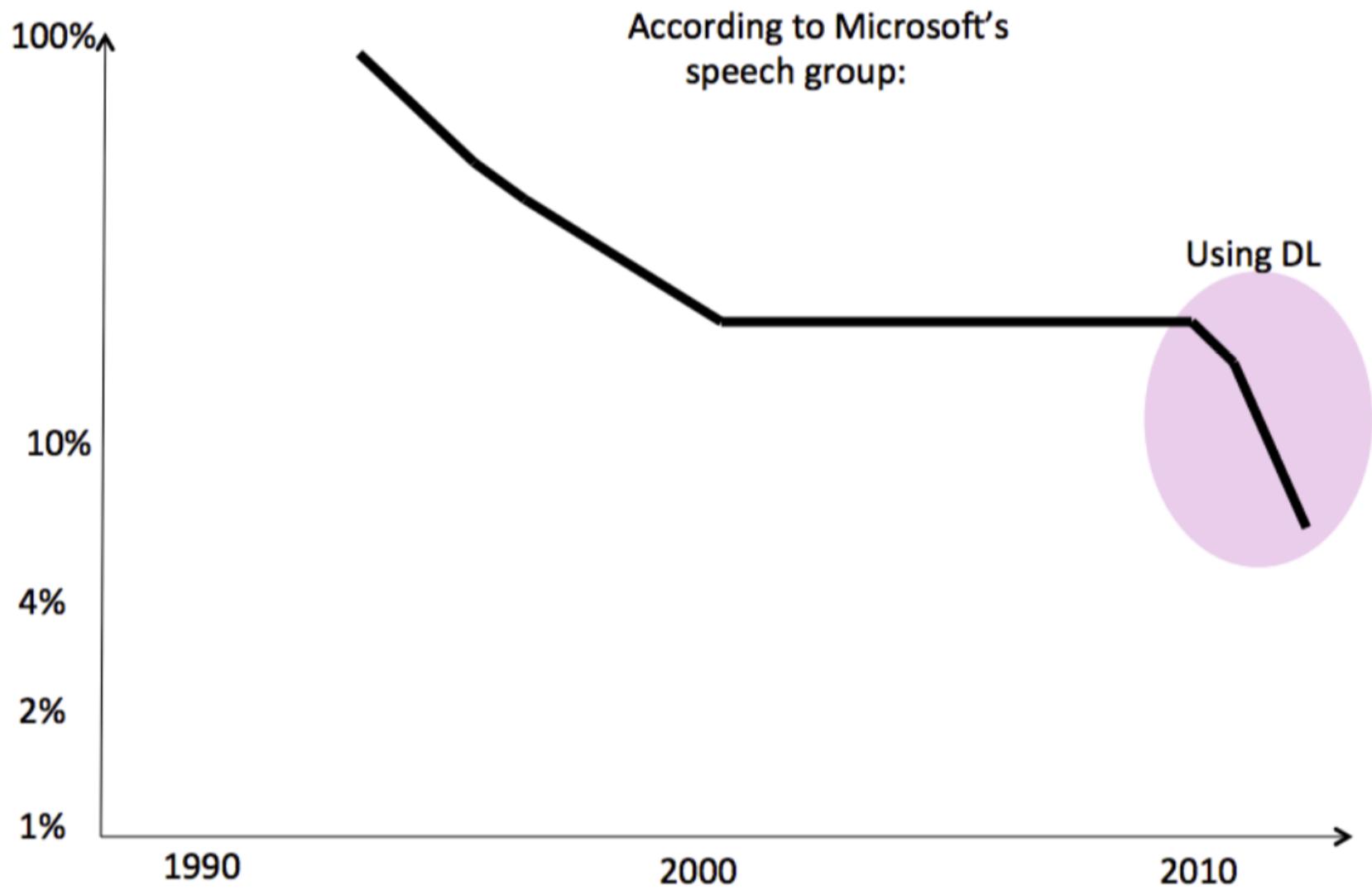
Deformable Part Model (2008)

Breakthroughs in 2012



ImageNet: The “computer vision World Cup”

Breakthroughs in 2012





www.image-net.org

22K categories and **14M** images

- Animals
 - Bird
 - Fish
 - Mammal
 - Invertebrate
- Plants
 - Tree
 - Flower
 - Food
 - Materials
- Structures
 - Artifact
 - Tools
 - Appliances
 - Structures
- Person
- Scenes
 - Indoor
 - Geological Formations
- Sport Activities

IMAGENET Large Scale Visual Recognition Challenge

Steel drum

The Image Classification Challenge:

1,000 object classes

1,431,167 images



Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle



Output:
Scale
T-shirt
Giant panda
Drumstick
Mud turtle

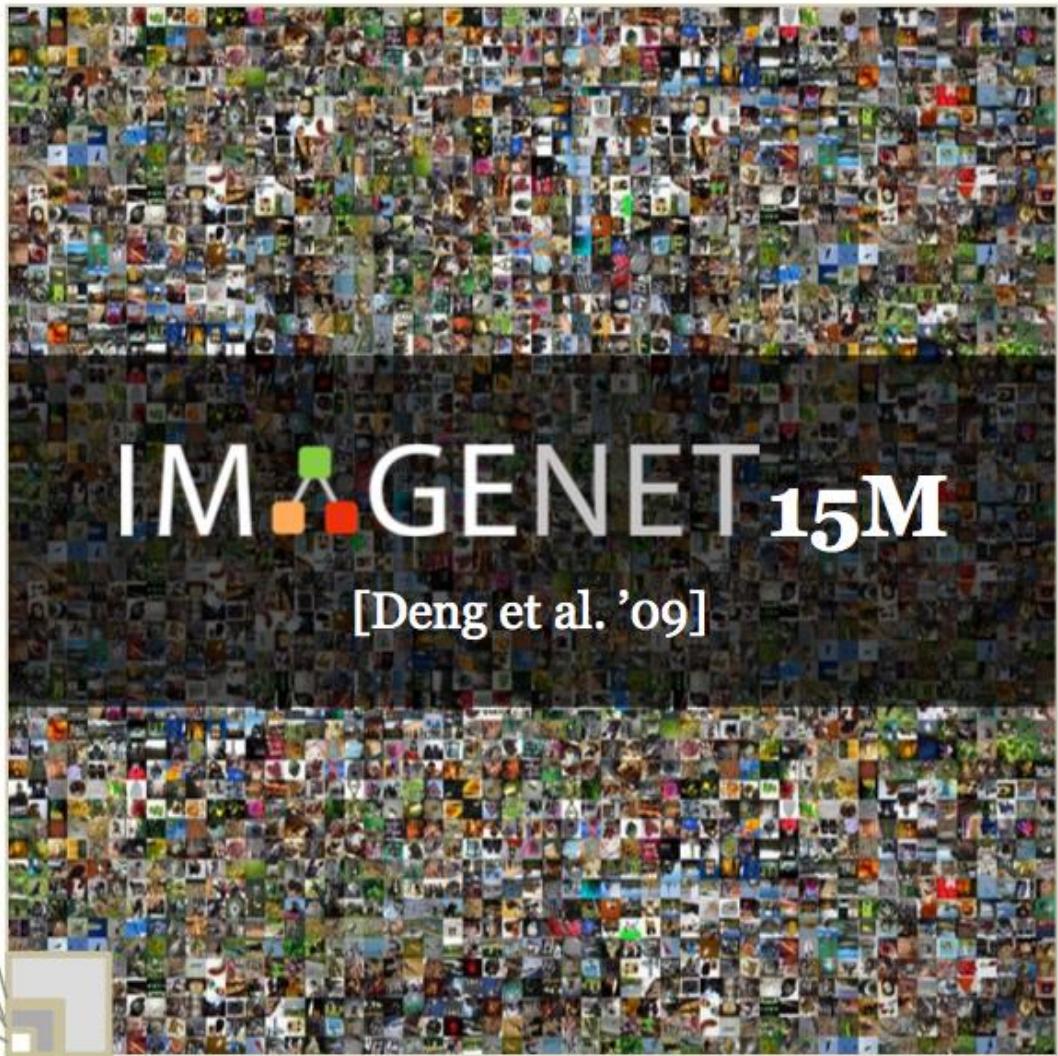


SUN, 131K
[Xiao et al. '10]

LabelMe, 37K
[Russell et al. '07]

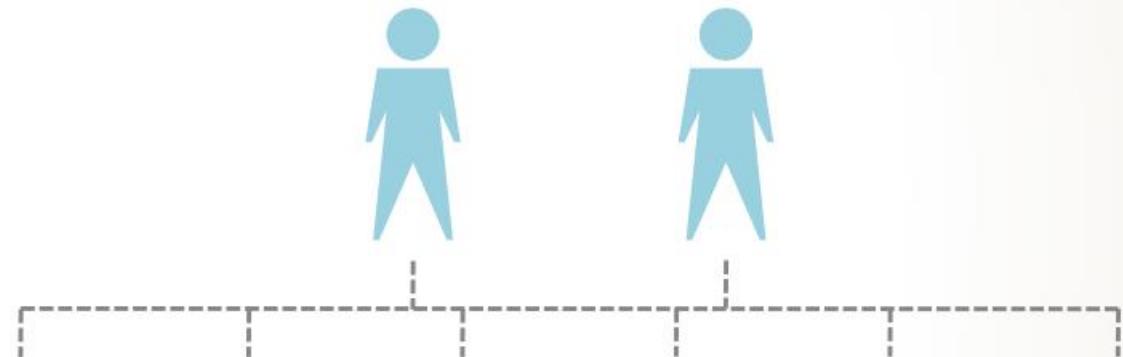
PASCAL VOC, 30K
[Everingham et al. '06-'12]

Caltech101, 9K
[Fei-Fei, Fergus, Perona, '03]



ImageNet PhD Students

**Crowdsourced
Labor**

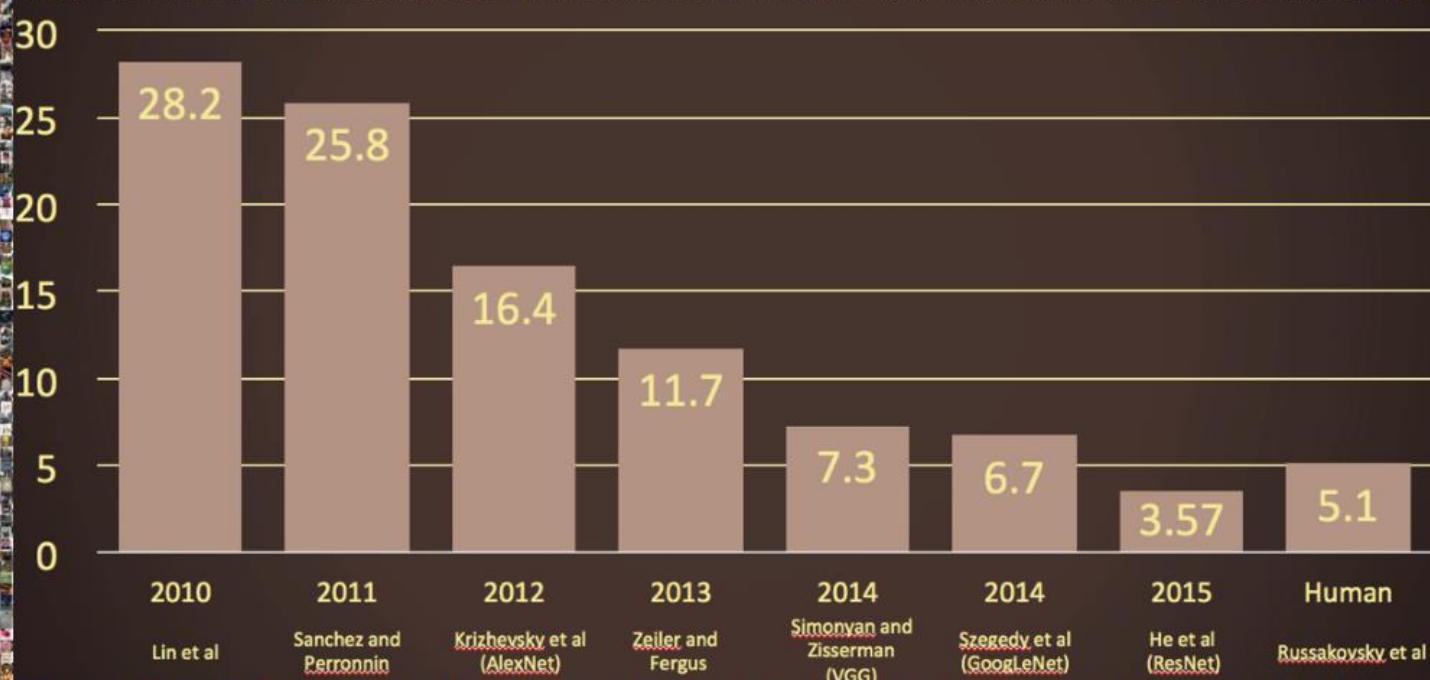


amazon mechanical turk™
Artificial Artificial Intelligence

**49k Workers *from 167 Countries*
2007-2010**

IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images



How Humans Compare

Human

5.1%

Top-5 error rate

GoogLeNet

6.8%

Top-5 error rate

Susceptible to:

- Fine-grained recognition
- Class unawareness
- Insufficient training data

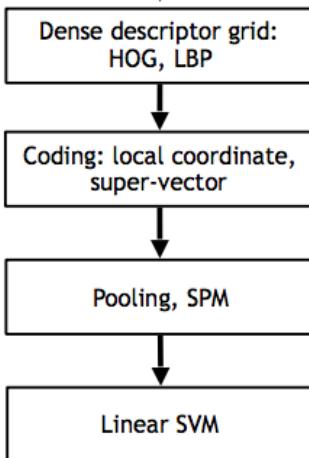
Susceptible to:

- Small, thin objects
- Image filters
- Abstract representations
- Miscellaneous sources

IMAGENET Large Scale Visual Recognition Challenge

Year 2010

NEC-UIUC

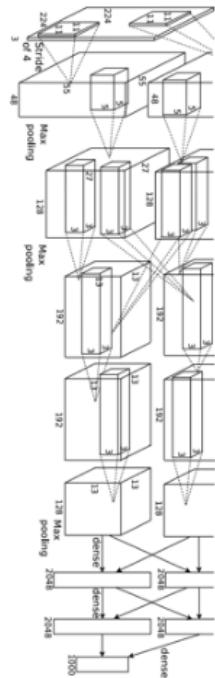


[Lin CVPR 2011]

Lion image by Swissfrog is licensed under CC BY 3.0

Year 2012

SuperVision



[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

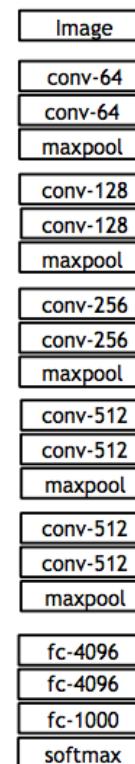
Year 2014

GoogLeNet



[Szegedy arxiv 2014]

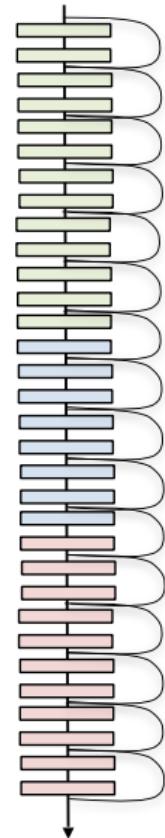
VGG



[Simonyan arxiv 2014]

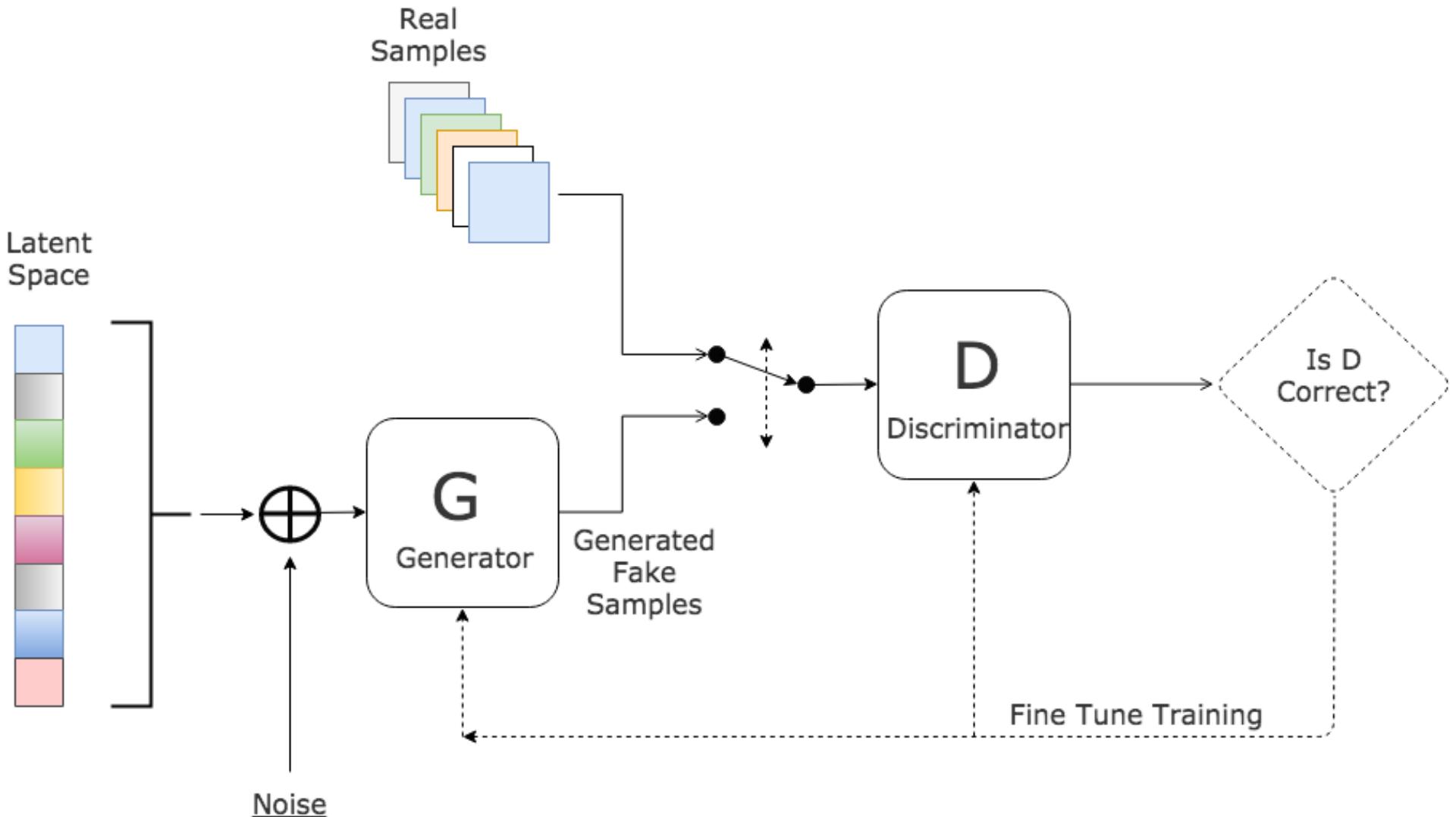
Year 2015

MSRA

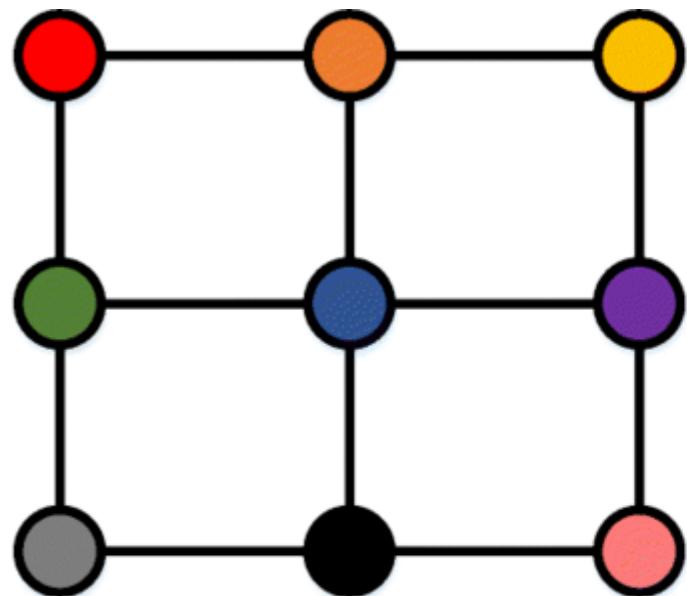


[He ICCV 2015]

Generative Adversarial Network (Ian Goodfellow, 2014)

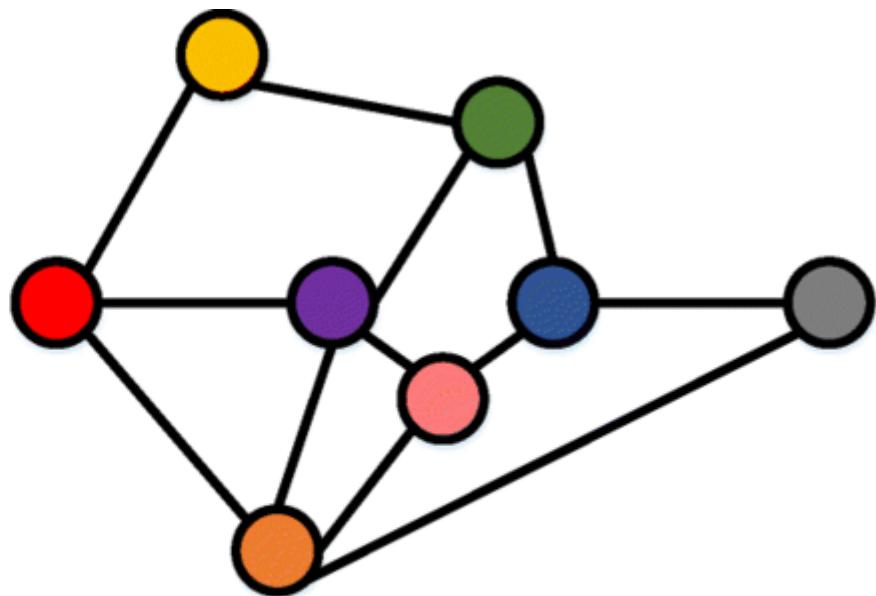


Graph Neural Networks



CNN

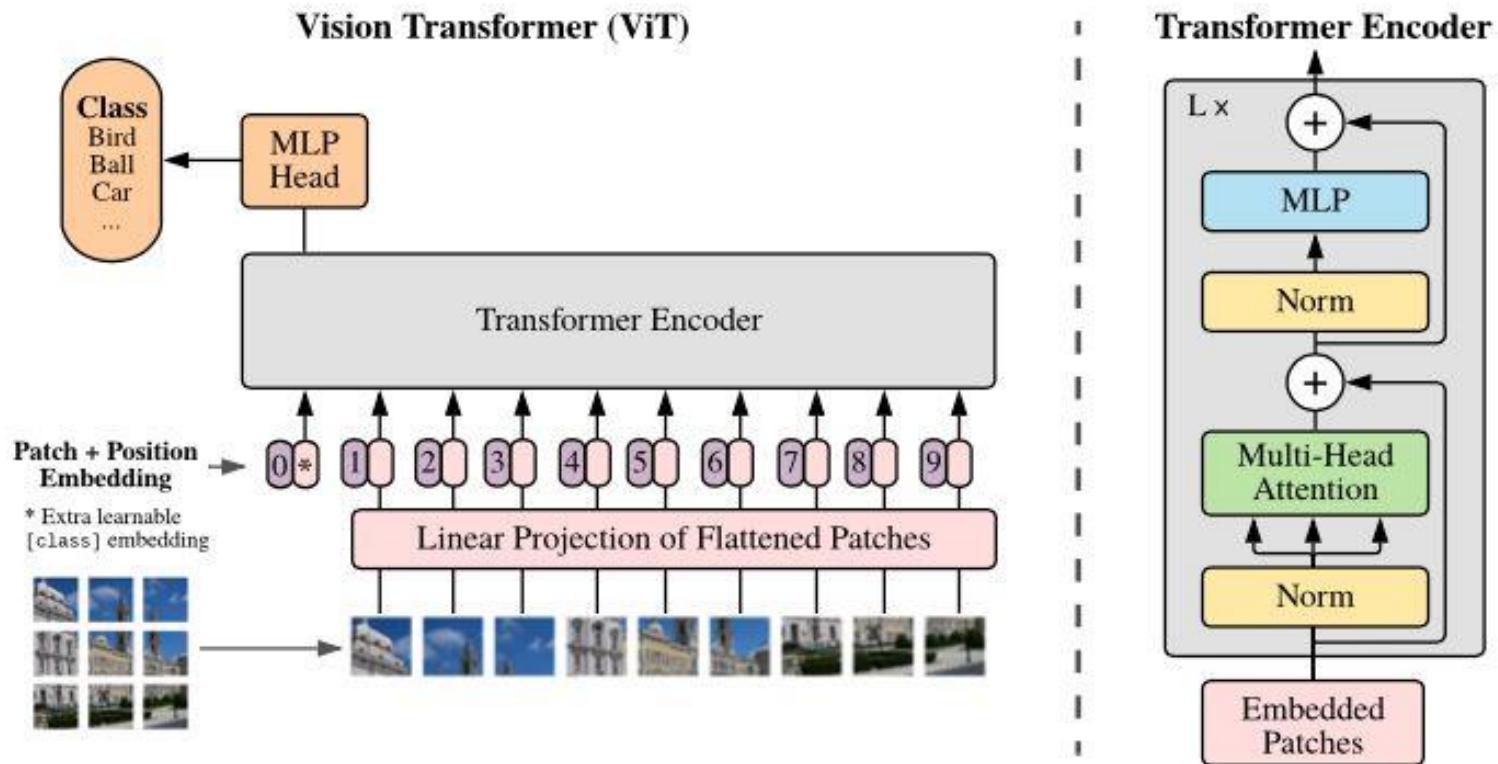
In Euclidean Space



GNN

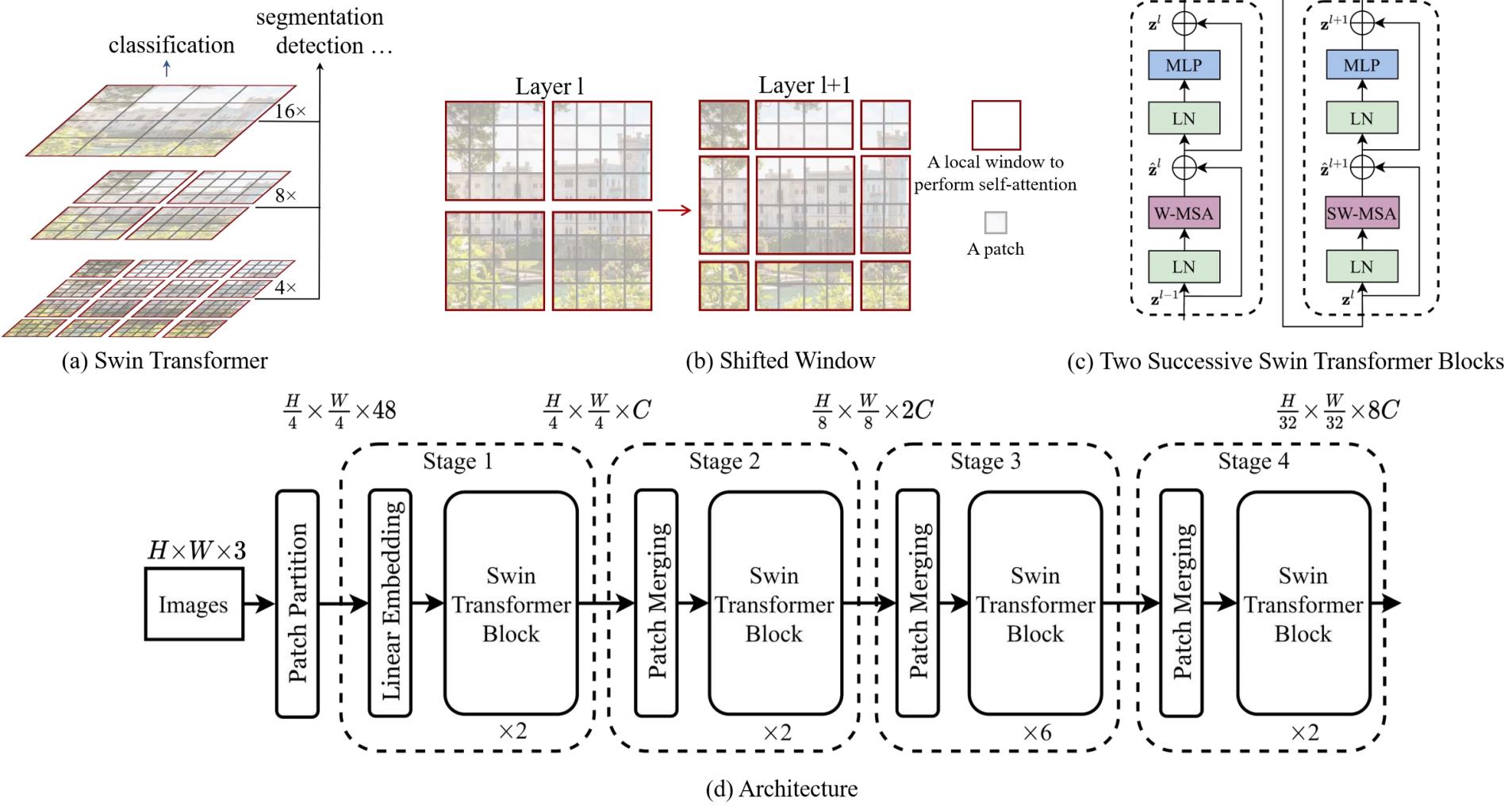
In Non-Euclidean Space

Vision Transformer (A. Dosovitskiy, 2020)

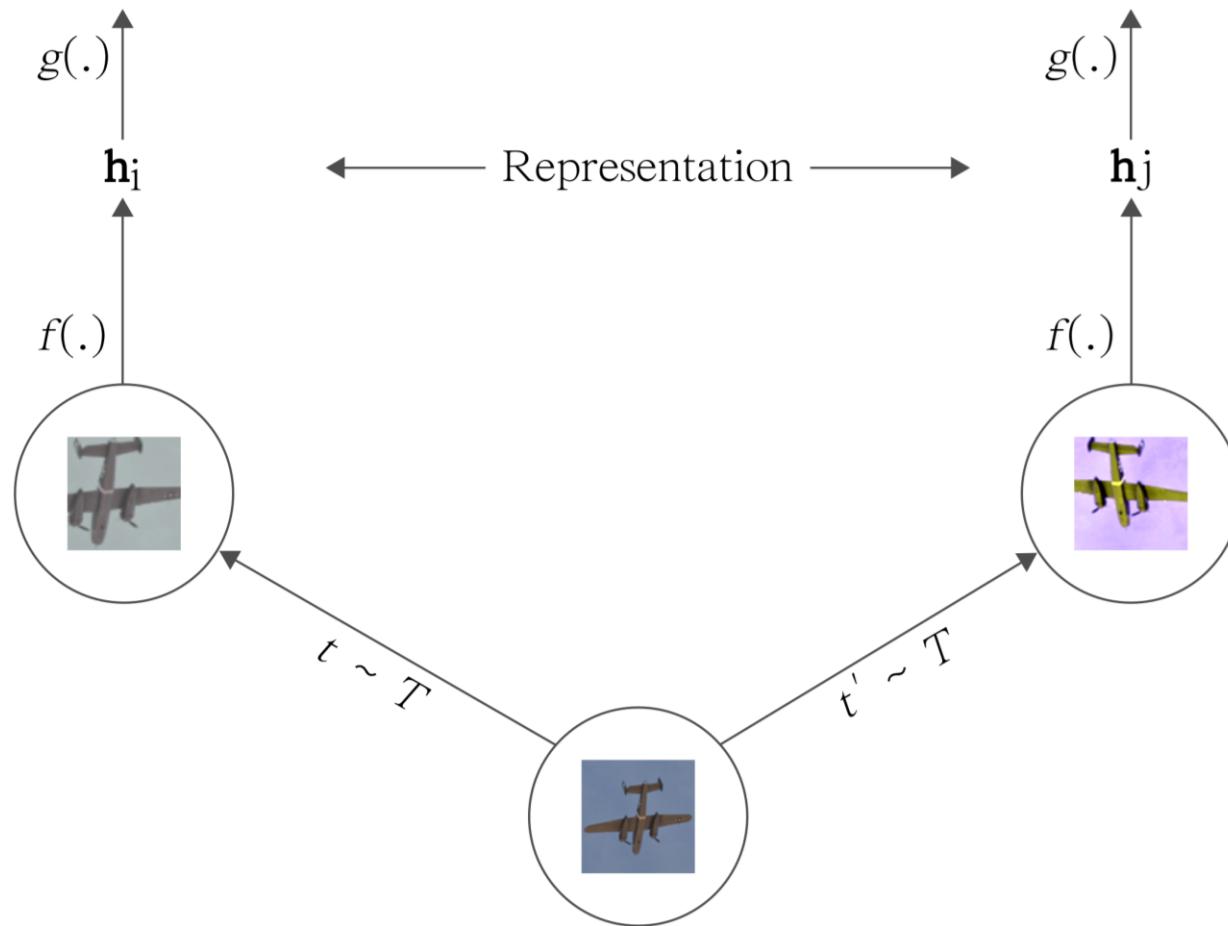


Hierarchical Vision Transformer

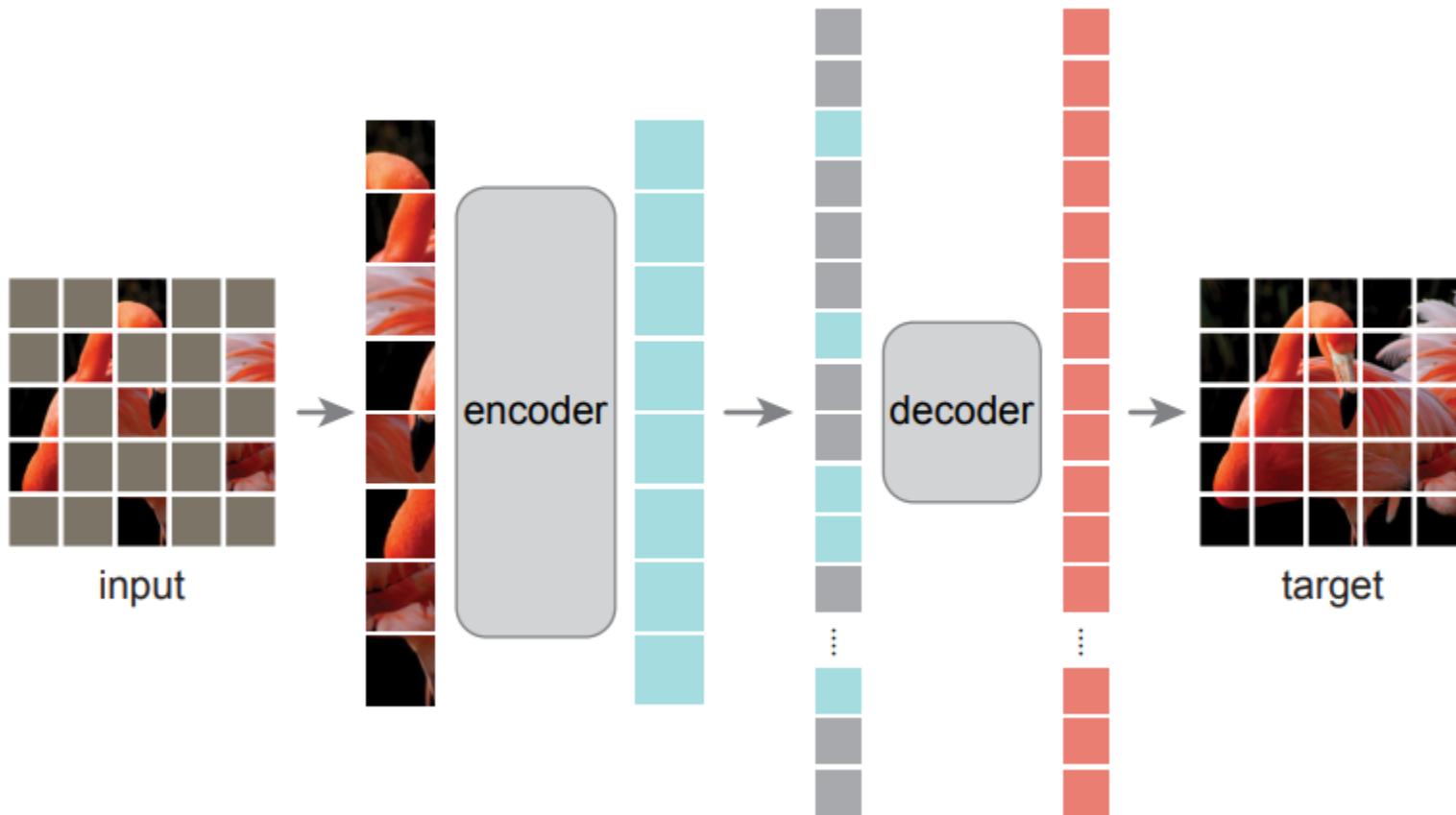
(Z. Liu, 2021)



Self-supervised Learning (Simclr, T. Chen, 2021)



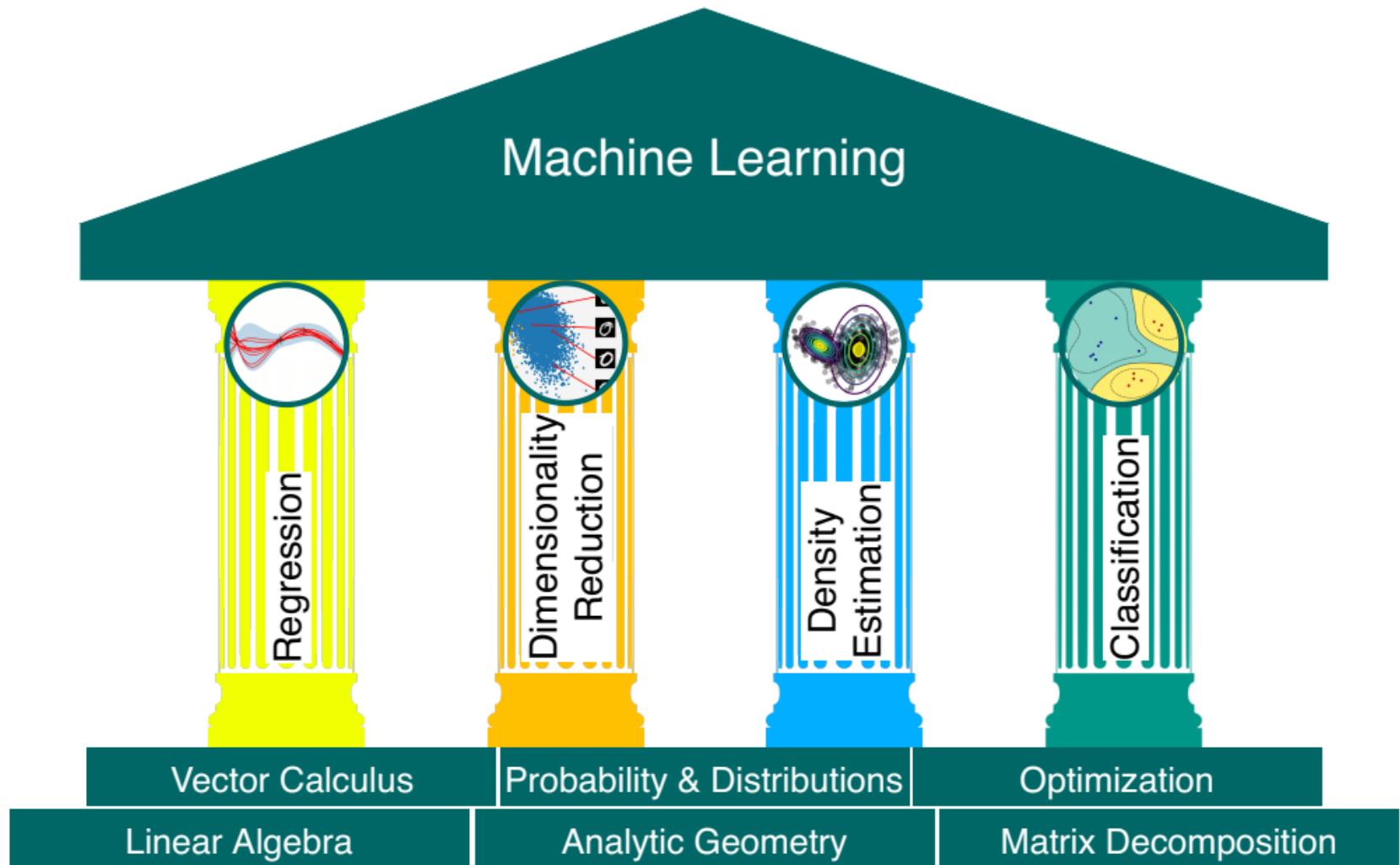
Self-supervised Learning (Masked Autoencoders, K. He, 2021)



Lecture One

Mathematics for AI

The foundations of Machine Learning



Main concept of Machine Learning

- We represent data as **vectors**.
- We choose an appropriate model, either using the **probabilistic or optimization** view.
- We learn from available data by using **numerical optimization** methods with the aim that the model performs well on data not used for training.

Essential Mathematical Concepts

Linear Algebra

We represent numerical data as vectors and represent a table of such data as a matrix. The study of vectors and matrices is called *linear algebra*.

■ Vectors

definition, scalars, addition, scalar multiplication, inner product(dot product), vector projection, cosine similarity, orthogonal vectors, normal and orthonormal vectors, vector norm, vector space, linear combination, linear span, linear independence, basis vectors

■ Matrices

definition, addition, transpose, scalar multiplication, matrix multiplication, matrix multiplication properties, hadamard product, functions, linear transformation, determinant, identity matrix, invertible matrix and inverse, rank, trace, popular type of matrices-symmetric, diagonal, orthogonal, orthonormal, positive definite matrix

■ Eigenvalues & eigenvectors

concept, intuition, significance, how to find

■ Principle component analysis

concept, properties, applications

■ Singular value decomposition

concept, properties, applications

Essential Mathematical Concepts

Calculus and Optimization

To train machine learning models, we typically find parameters that maximize some performance measure. Many optimization techniques require the concept of a gradient, which tells us the direction in which to search for a solution. Vector calculus is about gradient computation, and optimization is to find maxima/minima of functions.

- **Functions**
- **Scalar derivative**
definition, intuition, common rules of differentiation, chain rule, partial derivatives
- **Gradient**
concept, intuition, properties, directional derivative
- **Vector and matrix calculus**
how to find derivative of {scalar-valued, vector-valued} function wrt a {scalar, vector}
-> four combinations (Jacobian)
- **Gradient algorithms**
local/global maxima and minima, saddle point, convex functions, gradient descent algorithms (batch, mini-batch, stochastic, their performance comparison)

Essential Mathematical Concepts

Probability

We often would also like to have predictors that allow us to express some sort of uncertainty, e.g., to quantify the confidence we have about the value of the prediction at a particular test data point. Quantification of uncertainty is the realm of *probability theory*.

- **Basic rules and axioms**
events, sample space, frequentist approach, dependent and independent events, conditional probability
- **Random variables**
continuous and discrete, expectation, variance, distributions- joint and conditional
- **Bayes' Theorem, MAP, MLE**
- **Popular distributions**
binomial, bernoulli, poisson, exponential, gaussian
- **Conjugate priors**

Notation

$a, b, c, \alpha, \beta, \gamma$	Scalars are lowercase
$\mathbf{x}, \mathbf{y}, \mathbf{z}$	Vectors are bold lowercase
$\mathbf{A}, \mathbf{B}, \mathbf{C}$	Matrices are bold uppercase
$\mathbf{x}^\top, \mathbf{A}^\top$	Transpose of a vector or matrix
\mathbf{A}^{-1}	Inverse of a matrix
$\langle \mathbf{x}, \mathbf{y} \rangle$	Inner product of \mathbf{x} and \mathbf{y}
$\mathbf{x}^\top \mathbf{y}$	Dot product of \mathbf{x} and \mathbf{y}
$B = (\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3)$	(Ordered) tuple
$\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3]$	Matrix of column vectors stacked horizontally
$\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3\}$	Set of vectors (unordered)
\mathbb{Z}, \mathbb{N}	Integers and natural numbers, respectively
\mathbb{R}, \mathbb{C}	Real and complex numbers, respectively
\mathbb{R}^n	n -dimensional vector space of real numbers

Notation

$\forall x$	Universal quantifier: for all x
$\exists x$	Existential quantifier: there exists x
$a := b$	a is defined as b
$a =: b$	b is defined as a
$a \propto b$	a is proportional to b , i.e., $a = \text{constant} \cdot b$
$g \circ f$	Function composition: “ g after f ”
\iff	If and only if
\implies	Implies
\mathcal{A}, \mathcal{C}	Sets
$a \in \mathcal{A}$	a is an element of set \mathcal{A}
\emptyset	Empty set
$\mathcal{A} \setminus \mathcal{B}$	\mathcal{A} without \mathcal{B} : the set of elements in \mathcal{A} but not in \mathcal{B}
D	Number of dimensions; indexed by $d = 1, \dots, D$
N	Number of data points; indexed by $n = 1, \dots, N$

Notation

I_m	Identity matrix of size $m \times m$
$0_{m,n}$	Matrix of zeros of size $m \times n$
$1_{m,n}$	Matrix of ones of size $m \times n$
e_i	Standard/canonical vector (where i is the component that is 1)
dim	Dimensionality of vector space
$\text{rk}(\mathbf{A})$	Rank of matrix \mathbf{A}
$\text{Im}(\Phi)$	Image of linear mapping Φ
$\ker(\Phi)$	Kernel (null space) of a linear mapping Φ
$\text{span}[\mathbf{b}_1]$	Span (generating set) of \mathbf{b}_1
$\text{tr}(\mathbf{A})$	Trace of \mathbf{A}
$\det(\mathbf{A})$	Determinant of \mathbf{A}
$ \cdot $	Absolute value or determinant (depending on context)
$\ \cdot\ $	Norm; Euclidean, unless specified
λ	Eigenvalue or Lagrange multiplier
E_λ	Eigenspace corresponding to eigenvalue λ

Notation

$x \perp y$	Vectors x and y are orthogonal
V	Vector space
V^\perp	Orthogonal complement of vector space V
$\sum_{n=1}^N x_n$	Sum of the x_n : $x_1 + \dots + x_N$
$\prod_{n=1}^N x_n$	Product of the x_n : $x_1 \cdot \dots \cdot x_N$
θ	Parameter vector
$\frac{\partial f}{\partial x}$	Partial derivative of f with respect to x
$\frac{df}{dx}$	Total derivative of f with respect to x
∇	Gradient
$f_* = \min_x f(x)$	The smallest function value of f
$x_* \in \arg \min_x f(x)$	The value x_* that minimizes f (note: arg min returns a set of values)

Notation

\mathcal{L}	Lagrangian
\mathcal{L}	Negative log-likelihood
$\binom{n}{k}$	Binomial coefficient, n choose k
$\mathbb{V}_X[x]$	Variance of x with respect to the random variable X
$\mathbb{E}_X[x]$	Expectation of x with respect to the random variable X
$\text{Cov}_{X,Y}[x, y]$	Covariance between x and y .
$X \perp\!\!\!\perp Y Z$	X is conditionally independent of Y given Z
$X \sim p$	Random variable X is distributed according to p
$\mathcal{N}(\mu, \Sigma)$	Gaussian distribution with mean μ and covariance Σ
$\text{Ber}(\mu)$	Bernoulli distribution with parameter μ
$\text{Bin}(N, \mu)$	Binomial distribution with parameters N, μ
$\text{Beta}(\alpha, \beta)$	Beta distribution with parameters α, β

Notation

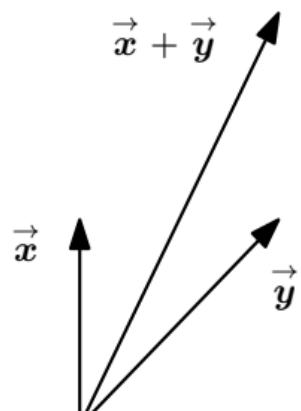
\mathcal{L}	Lagrangian
\mathcal{L}	Negative log-likelihood
$\binom{n}{k}$	Binomial coefficient, n choose k
$\mathbb{V}_X[x]$	Variance of x with respect to the random variable X
$\mathbb{E}_X[x]$	Expectation of x with respect to the random variable X
$\text{Cov}_{X,Y}[x, y]$	Covariance between x and y .
$X \perp\!\!\!\perp Y Z$	X is conditionally independent of Y given Z
$X \sim p$	Random variable X is distributed according to p
$\mathcal{N}(\mu, \Sigma)$	Gaussian distribution with mean μ and covariance Σ
$\text{Ber}(\mu)$	Bernoulli distribution with parameter μ
$\text{Bin}(N, \mu)$	Binomial distribution with parameters N, μ
$\text{Beta}(\alpha, \beta)$	Beta distribution with parameters α, β

Linear Algebra

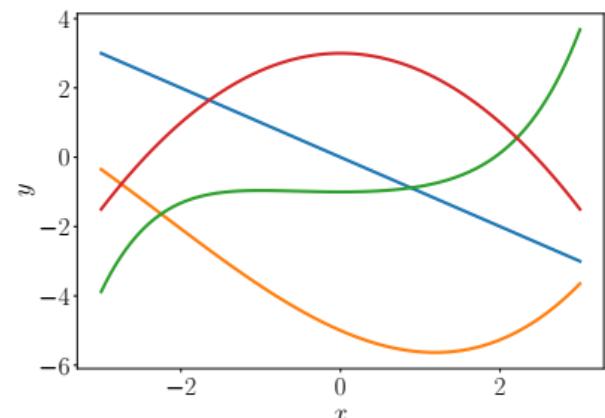
Vectors

In general, vectors are special objects that can be added together and multiplied by scalars to produce another object of the same kind. From an abstract mathematical viewpoint, any object that satisfies these two properties can be considered a vector.

1. Geometric vectors
2. Polynomials
3. Audio signals
4. Elements of \mathbb{R}^n



(a) Geometric vectors.



(b) Polynomials.

Matrices: definition

Definition 2.1 (Matrix). With $m, n \in \mathbb{N}$ a real-valued (m, n) matrix A is an $m \cdot n$ -tuple of elements a_{ij} , $i = 1, \dots, m$, $j = 1, \dots, n$, which is ordered according to a rectangular scheme consisting of m rows and n columns:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad a_{ij} \in \mathbb{R}. \quad (2.11)$$

By convention $(1, n)$ -matrices are called *rows* and $(m, 1)$ -matrices are called *columns*. These special matrices are also called *row/column vectors*.

Matrices: addition and multiplication

The sum of two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m \times n}$ is defined as the element-wise sum, i.e.,

$$\mathbf{A} + \mathbf{B} := \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}. \quad (2.12)$$

For matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times k}$, the elements c_{ij} of the product $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times k}$ are computed as

$$c_{ij} = \sum_{l=1}^n a_{il}b_{lj}, \quad i = 1, \dots, m, \quad j = 1, \dots, k. \quad (2.13)$$

Q: What is the computational complexity of matrix addition and multiplication?

Matrices: addition and multiplication

The sum of two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m \times n}$ is defined as the element-wise sum, i.e.,

$$\mathbf{A} + \mathbf{B} := \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}. \quad (2.12)$$

For matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times k}$, the elements c_{ij} of the product $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times k}$ are computed as

$$c_{ij} = \sum_{l=1}^n a_{il}b_{lj}, \quad i = 1, \dots, m, \quad j = 1, \dots, k. \quad (2.13)$$

Q: What is the computational complexity of matrix addition and multiplication? O(mn) and O(mnk)

Matrices: addition and multiplication

Remark1: Matrices can only be multiplied if their “neighboring” dimensions match.

Remark2: Matrix multiplication is not defined as an element-wise operation on matrix elements. This kind of element-wise multiplication is called a [Hadamard product](#).

Remark3: Even if both matrix multiplications \mathbf{AB} and \mathbf{BA} are defined, the dimensions of the results can be different.

$$\begin{array}{c|c} \begin{matrix} \text{blue} & \text{blue} \\ \text{blue} & \text{blue} \\ \text{blue} & \text{blue} \end{matrix} & \begin{matrix} \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} \end{matrix} \\ \hline & = \\ \hline \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} & \end{array} \quad \begin{array}{c|c} \begin{matrix} \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} \end{matrix} & \begin{matrix} \text{blue} & \text{blue} \\ \text{blue} & \text{blue} \\ \text{blue} & \text{blue} \end{matrix} \\ \hline & = \\ \hline \begin{matrix} \text{green} & \text{green} \\ \text{green} & \text{green} \end{matrix} & \end{array}$$

Matrices: addition and multiplication

- *Associativity:*

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p}, \mathbf{C} \in \mathbb{R}^{p \times q} : (\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC}) \quad (2.18)$$

- *Distributivity:*

$$\forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{n \times p} : (\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC} \quad (2.19a)$$

$$\mathbf{A}(\mathbf{C} + \mathbf{D}) = \mathbf{AC} + \mathbf{AD} \quad (2.19b)$$

- Multiplication with the identity matrix:

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n} : \mathbf{I}_m \mathbf{A} = \mathbf{A} \mathbf{I}_n = \mathbf{A} \quad (2.20)$$

Note that $\mathbf{I}_m \neq \mathbf{I}_n$ for $m \neq n$.

Matrices: identity matrix

Definition 2.2 (Identity Matrix). In $\mathbb{R}^{n \times n}$, we define the *identity matrix*

$$\mathbf{I}_n := \begin{bmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (2.17)$$

Matrices: inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $A \in \mathbb{R}^{n \times n}$. Let matrix $B \in \mathbb{R}^{n \times n}$ have the property that $AB = I_n = BA$. B is called the *inverse* of A and denoted by A^{-1} .

Matrices: inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $A \in \mathbb{R}^{n \times n}$. Let matrix $B \in \mathbb{R}^{n \times n}$ have the property that $AB = I_n = BA$. B is called the *inverse* of A and denoted by A^{-1} .

Q: Is every matrix possesses an inverse?

Matrices: inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $A \in \mathbb{R}^{n \times n}$. Let matrix $B \in \mathbb{R}^{n \times n}$ have the property that $AB = I_n = BA$. B is called the *inverse* of A and denoted by A^{-1} .

Q: Is every matrix possesses an inverse? No. Positive Definite Matrices are always invertible.

Matrices: inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. Let matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ have the property that $\mathbf{AB} = \mathbf{I}_n = \mathbf{BA}$. \mathbf{B} is called the *inverse* of \mathbf{A} and denoted by \mathbf{A}^{-1} .

Definition 2.4 (Transpose). For $\mathbf{A} \in \mathbb{R}^{m \times n}$ the matrix $\mathbf{B} \in \mathbb{R}^{n \times m}$ with $b_{ij} = a_{ji}$ is called the *transpose* of \mathbf{A} . We write $\mathbf{B} = \mathbf{A}^\top$.

In general, \mathbf{A}^\top can be obtained by writing the columns of \mathbf{A} as the rows of \mathbf{A}^\top . The following are important properties of inverses and transposes:

$$\mathbf{AA}^{-1} = \mathbf{I} = \mathbf{A}^{-1}\mathbf{A} \quad (2.26)$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1} \quad (2.27)$$

$$(\mathbf{A} + \mathbf{B})^{-1} \neq \mathbf{A}^{-1} + \mathbf{B}^{-1} \quad (2.28)$$

$$(\mathbf{A}^\top)^\top = \mathbf{A} \quad (2.29)$$

$$(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top \quad (2.30)$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top \quad (2.31)$$

Matrices: multiplication by a scalar

Let us look at what happens to matrices when they are multiplied by a scalar $\lambda \in \mathbb{R}$. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\lambda \in \mathbb{R}$. Then $\lambda\mathbf{A} = \mathbf{K}$, $K_{ij} = \lambda a_{ij}$. Practically, λ scales each element of \mathbf{A} . For $\lambda, \psi \in \mathbb{R}$, the following holds:

- *Associativity:*

$$(\lambda\psi)\mathbf{C} = \lambda(\psi\mathbf{C}), \quad \mathbf{C} \in \mathbb{R}^{m \times n}$$

- $\lambda(\mathbf{B}\mathbf{C}) = (\lambda\mathbf{B})\mathbf{C} = \mathbf{B}(\lambda\mathbf{C}) = (\mathbf{B}\mathbf{C})\lambda, \quad \mathbf{B} \in \mathbb{R}^{m \times n}, \mathbf{C} \in \mathbb{R}^{n \times k}$.

Note that this allows us to move scalar values around.

- $(\lambda\mathbf{C})^\top = \mathbf{C}^\top\lambda^\top = \mathbf{C}^\top\lambda = \lambda\mathbf{C}^\top$ since $\lambda = \lambda^\top$ for all $\lambda \in \mathbb{R}$.

- *Distributivity:*

$$(\lambda + \psi)\mathbf{C} = \lambda\mathbf{C} + \psi\mathbf{C}, \quad \mathbf{C} \in \mathbb{R}^{m \times n}$$

$$\lambda(\mathbf{B} + \mathbf{C}) = \lambda\mathbf{B} + \lambda\mathbf{C}, \quad \mathbf{B}, \mathbf{C} \in \mathbb{R}^{m \times n}$$

Matrix of linear transformation

3d linear transformations

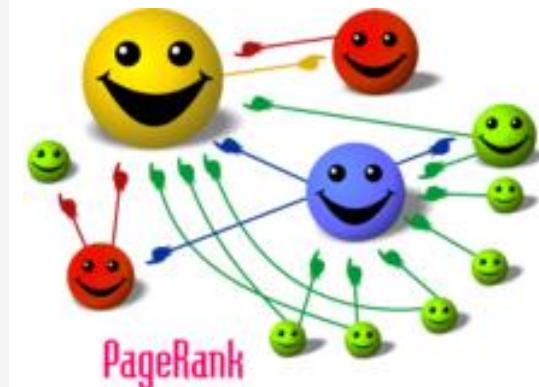
Matrix multiplication as composition

Eigenvectors and Eigenvalues



Google's PageRank

Google uses the eigenvector corresponding to the maximal eigenvalue of a matrix \mathbf{A} to determine the rank of a page for search. The idea for the PageRank algorithm, developed at Stanford University by Larry Page and Sergey Brin in 1996, was that the importance of any web page can be approximated by the importance of pages that link to it. For this, they write down all web sites as a huge directed graph that shows which page links to which. PageRank computes the weight (importance) $x_i \geq 0$ of a web site a_i by counting the number of pages pointing to a_i . Moreover, PageRank takes into account the importance of the web sites that link to a_i . The navigation behavior of a user is then modeled by a transition matrix \mathbf{A} of this graph that tells us with what (click) probability somebody will end up on a different web site. The matrix \mathbf{A} has the property that for any initial rank/importance vector \mathbf{x} of a web site the sequence $\mathbf{x}, \mathbf{Ax}, \mathbf{A}^2\mathbf{x}, \dots$ converges to a vector \mathbf{x}^* . This vector is called the *PageRank* and satisfies $\mathbf{Ax}^* = \mathbf{x}^*$, i.e., it is an eigenvector (with corresponding eigenvalue 1) of \mathbf{A} . After normalizing \mathbf{x}^* , such that $\|\mathbf{x}^*\| = 1$, we can interpret the entries as probabilities. More details and different perspectives on PageRank can be found in the original technical report (Page et al., 1999).



Low-rank matrix approximation

We considered the SVD as a way to factorize $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top \in \mathbb{R}^{m \times n}$ into the product of three matrices, where $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ are orthogonal and Σ contains the singular values on its main diagonal. Instead of doing the full SVD factorization, we will now investigate how the SVD allows us to represent a matrix \mathbf{A} as a sum of simpler (low-rank) matrices \mathbf{A}_i , which lends itself to a matrix approximation scheme that is cheaper to compute than the full SVD.

We construct a rank-1 matrix $\mathbf{A}_i \in \mathbb{R}^{m \times n}$ as

$$\mathbf{A}_i := \mathbf{u}_i \mathbf{v}_i^\top, \quad (4.90)$$

which is formed by the outer product of the i th orthogonal column vector of \mathbf{U} and \mathbf{V} . Figure 4.11 shows an image of Stonehenge, which can be represented by a matrix $\mathbf{A} \in \mathbb{R}^{1432 \times 1910}$, and some outer products \mathbf{A}_i , as defined in (4.90).

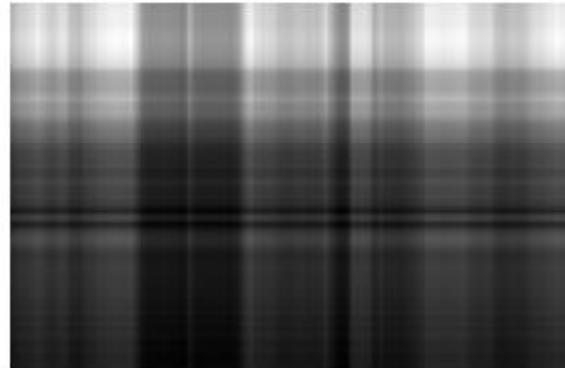
A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank r can be written as a sum of rank-1 matrices \mathbf{A}_i so that

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \sum_{i=1}^r \sigma_i \mathbf{A}_i, \quad (4.91)$$

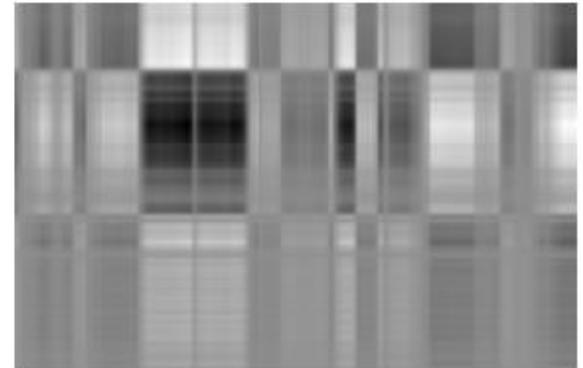
Low-rank matrix approximation



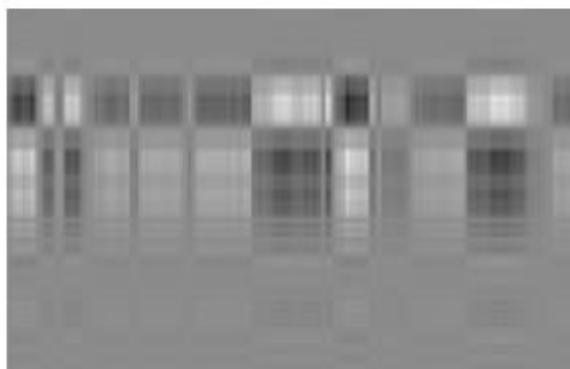
(a) Original image \mathbf{A} .



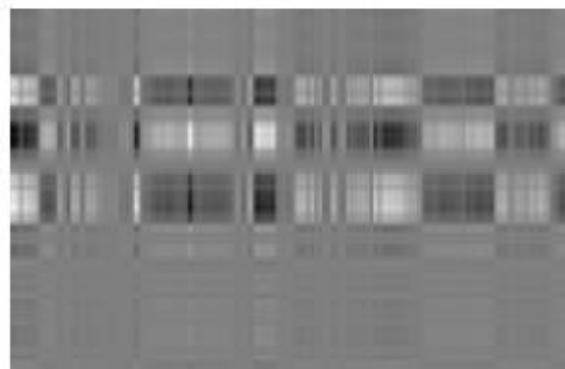
(b) \mathbf{A}_1 , $\sigma_1 \approx 228,052$.



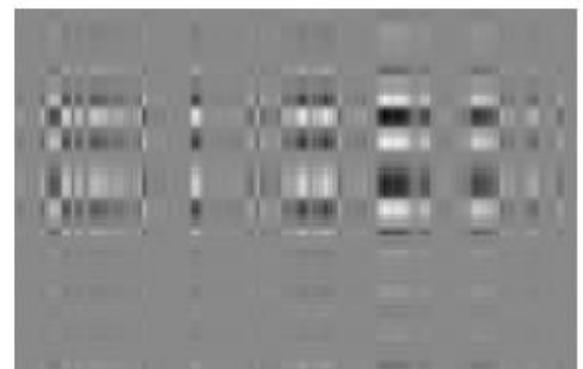
(c) \mathbf{A}_2 , $\sigma_2 \approx 40,647$.



(d) \mathbf{A}_3 , $\sigma_3 \approx 26,125$.



(e) \mathbf{A}_4 , $\sigma_4 \approx 20,232$.

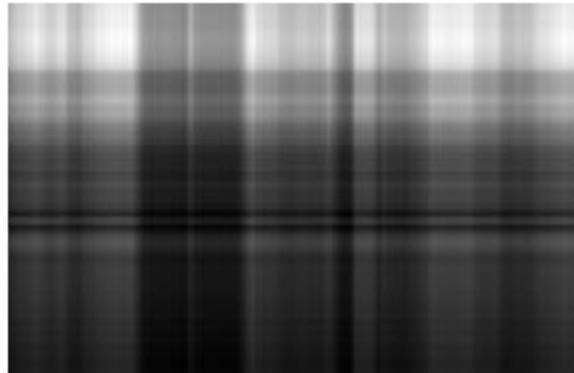


(f) \mathbf{A}_5 , $\sigma_5 \approx 15,436$.

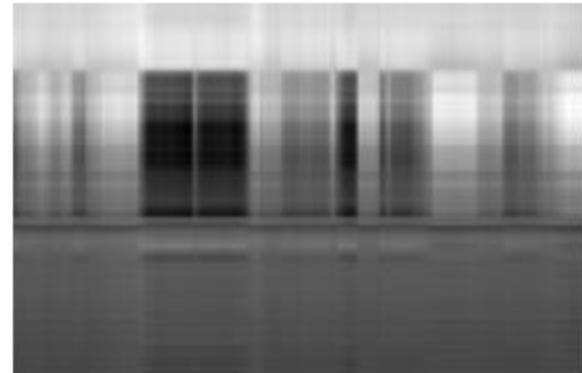
Low-rank matrix approximation



(a) Original image \mathbf{A} .



(b) Rank-1 approximation $\widehat{\mathbf{A}}(1)$.



(c) Rank-2 approximation $\widehat{\mathbf{A}}(2)$.



(d) Rank-3 approximation $\widehat{\mathbf{A}}(3)$.



(e) Rank-4 approximation $\widehat{\mathbf{A}}(4)$.



(f) Rank-5 approximation $\widehat{\mathbf{A}}(5)$.

Vector Calculus

Differentiation Rules

In the following, we briefly state basic differentiation rules, where we denote the derivative of f by f' .

Product rule: $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$ (5.29)

Quotient rule: $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$ (5.30)

Sum rule: $(f(x) + g(x))' = f'(x) + g'(x)$ (5.31)

Chain rule: $(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$ (5.32)

Here, $g \circ f$ denotes function composition $x \mapsto f(x) \mapsto g(f(x))$.

Partial Differentiation and Gradients

Definition 5.5 (Partial Derivative). For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^n$ of n variables x_1, \dots, x_n we define the *partial derivatives* as

$$\begin{aligned}\frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(\mathbf{x})}{h}\end{aligned}\tag{5.39}$$

and collect them in the row vector

$$\nabla_{\mathbf{x}} f = \text{grad } f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}, \tag{5.40}$$

Gradients of Vector-Valued Functions

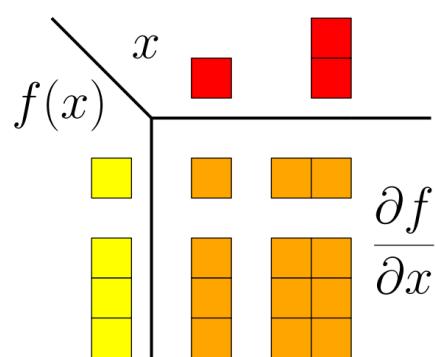
Definition 5.6 (Jacobian). The collection of all first-order partial derivatives of a vector-valued function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called the *Jacobian*. The Jacobian \mathbf{J} is an $m \times n$ matrix, which we define and arrange as follows:

$$\mathbf{J} = \nabla_{\mathbf{x}} \mathbf{f} = \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \end{bmatrix} \quad (5.57)$$

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}, \quad (5.58)$$

Figure 5.6
Dimensionality of
(partial) derivatives.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad J(i, j) = \frac{\partial f_i}{\partial x_j}. \quad (5.59)$$



Chain Rule for vector differentiation

向量求导的链式法则：

若多个向量的依赖关系为 $u \rightarrow v \rightarrow w$.

$$\text{则 } \frac{\partial w}{\partial u} = \frac{\partial w}{\partial v} \cdot \frac{\partial v}{\partial u}$$

$\downarrow \quad \downarrow \quad \downarrow$
 $m \times 1 \quad n \times 1 \quad l \times 1$

$$\downarrow \quad \downarrow \quad \downarrow$$

$l \times m \quad l \times n \quad n \times m$

证明：
$$\frac{\partial w_i}{\partial u_j} = \sum_k \frac{\partial w_i}{\partial v_k} \cdot \frac{\partial v_k}{\partial u_j} \leftarrow u_j \begin{matrix} \nearrow v_1 \\ \vdots \\ \nearrow v_n \end{matrix} \rightarrow w_i$$

该法则可推广到多层中间变量的情形

A simplest example

Example 5.9 (Gradient of a Vector-Valued Function)

We are given

$$\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}, \quad \mathbf{f}(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N.$$

To compute the gradient $d\mathbf{f}/d\mathbf{x}$ we first determine the dimension of $d\mathbf{f}/d\mathbf{x}$: Since $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$, it follows that $d\mathbf{f}/d\mathbf{x} \in \mathbb{R}^{M \times N}$. Second, to compute the gradient we determine the partial derivatives of f with respect to every x_j :

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \implies \frac{\partial f_i}{\partial x_j} = A_{ij} \quad (5.67)$$

We collect the partial derivatives in the Jacobian and obtain the gradient

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = \mathbf{A} \in \mathbb{R}^{M \times N}. \quad (5.68)$$

More cases

① $\nabla Ax = A$, 特別地, $\nabla x = I$

证: $y = Ax \Rightarrow \frac{\partial y_i}{\partial x_j} = \frac{\partial \sum A_{ij}x_j}{\partial x_j} = A_{ij}$

② $\nabla(a^T x) = a^T$

证: $\frac{\partial a^T x}{\partial x_i} = \frac{\partial \sum a_j x_j}{\partial x_i} = \frac{\partial a_i x_i}{\partial x_i} = a_i$

③ $\nabla \|x\|_2^2 = \nabla(x^T x) = 2x$

证: $\frac{\partial \|x\|_2^2}{\partial x_i} = \frac{\partial \sum x_j^2}{\partial x_i} = \frac{\partial x_i^2}{\partial x_i} = 2x_i$

More cases

$$\textcircled{4} \quad \nabla(x^T A x) = x^T (A + A^T)$$

证:

$$f(x) = x^T A x$$

↓ 添加下标

$$f = x_1^T A x_2$$

↓ 分别求导

$$\frac{\partial f}{\partial x_1} = x_2^T A^T, \quad \frac{\partial f}{\partial x_2} = x_1^T A$$

↓ 去掉下标后相加

$$\frac{\partial f}{\partial x} = x^T (A + A^T)$$

More cases

$$⑤ \nabla(u^T v) = v^T (\nabla_x u) + u^T (\nabla_x v)$$

证 (利用变量多次出现的求导法则 (x 在 u, v 中出现).)

$$⑥ \nabla_x (\lambda(x) f(x)) = f(x) \cdot \nabla_x \lambda(x) + \lambda(x) \cdot \nabla_x f(x)$$

$$\text{证 } \left(\frac{\partial \lambda f_i}{\partial x_j} = f_i \cdot \frac{\partial \lambda}{\partial x_j} + \lambda \cdot \frac{\partial f_i}{\partial x_j} \right).$$

More cases

$$\textcircled{1} \quad \nabla \text{tr}(A^T X) = \nabla \text{tr}(AX^T) = A$$

证

$$\frac{\partial \text{tr}(A^T X)}{\partial X_{ij}} = \frac{\partial \sum_{i,j} (A_{ij} \cdot X_{ij})}{\partial X_{ij}} = A_{ij}$$

$$\textcircled{2} \quad \nabla \text{tr}(XAX^T B) = B^T X A^T + B X A$$

证 (利用变量多次出现的求导法则)

$$\textcircled{3} \quad \nabla a^T X b = ab^T$$

证

$$\text{LHS} = \nabla \text{tr}(a^T X b) = \nabla \text{tr}(X b a^T) = ab^T = \text{RHS}$$

More cases

⑩ $\nabla_X a^T X^T X a = 2Xaa^T$ (直接用⑨证明)

⑪ $\nabla_X (Xa-b)^T (Xa-b) = 2(Xa-b)a^T$ (展开后用⑨证明)

⑫ $\nabla_X \|XA^T - B\|_F^2 = 2(XA^T - B)A$

(注意 $\|A\|_F^2 = \sum_{i,j} A_{ij}^2 = \text{tr}(A^T A)$)

⑬ 行列式: $\nabla_X |X| = |X| (X^{-1})^T$

More cases

⑭ 设 $y = f(u)$, $u = g(X)$, 则 $\frac{\partial y}{\partial X} = \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial X}$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
参数 参数 $m \times n$ $|X|$ $m \times n$.

⑮ 设 $f(Y) : R^{m \times p} \rightarrow R$, 则 $\begin{cases} \nabla_X f(AX+B) = A^T \nabla_{Y=AX+B} f \\ \nabla_X f(XC+D) = (\nabla_{Y=XG+D} f) \cdot C^T \end{cases}$

$\downarrow \quad \downarrow \quad \downarrow$
 $m \times n \quad n \times p \quad m \times p$

⑯ $\nabla_X f(Ax+b) = A^T \nabla_{Y=Ax+b} f$

(⑮ 的向量版本)

Useful Tricks

Trick: turn all to vector/matrix operations.

$$\left\{ \begin{array}{l} \sum_i u_i v_i = u^T v \\ \sum_{ij} A_{ij} B_{ij} = \text{tr}(AB^T) \\ \sum_i x_i^2 = x^T x \\ \|A\|_F^2 = \text{tr}(AA^T) \end{array} \right.$$

Useful Identities

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^\top = \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right)^\top \quad (5.99)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{f}(\mathbf{X})) = \text{tr} \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.100)$$

$$\frac{\partial}{\partial \mathbf{X}} \det(\mathbf{f}(\mathbf{X})) = \det(\mathbf{f}(\mathbf{X})) \text{tr} \left(\mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.101)$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} = -\mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} \quad (5.102)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^\top \mathbf{a} \mathbf{b}^\top (\mathbf{X}^{-1})^\top \quad (5.103)$$

$$\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.104)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.105)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top \quad (5.106)$$

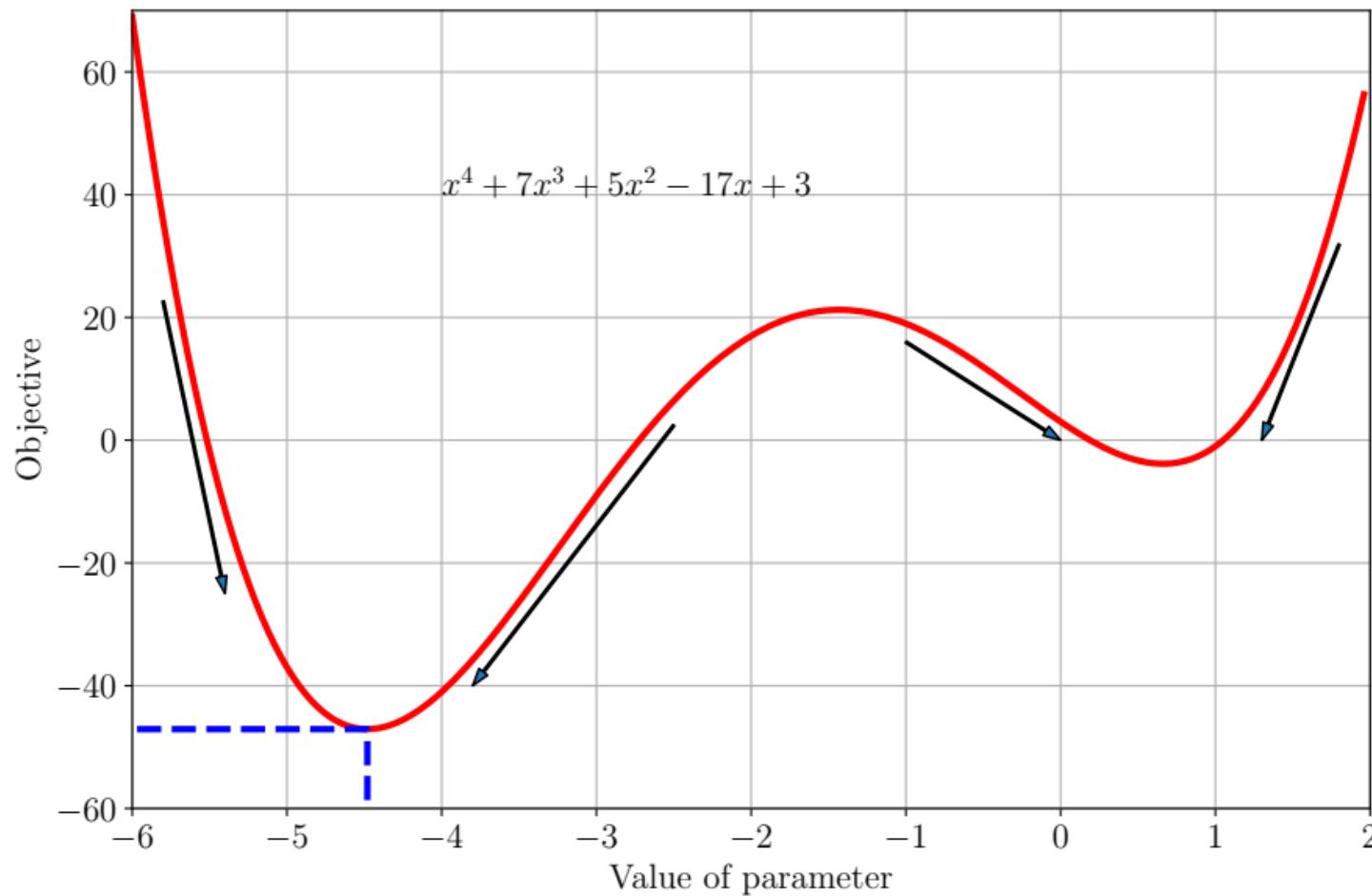
$$\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^\top (\mathbf{B} + \mathbf{B}^\top) \quad (5.107)$$

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W} (\mathbf{x} - \mathbf{A}\mathbf{s}) = -2(\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W} \mathbf{A} \quad \text{for symmetric } \mathbf{W}$$
$$(5.108)$$

Optimization

Machine Learning as Optimization

- The process of machine learning is to find parameters that minimize objective



Gradient Descent

We now consider the problem of solving for the minimum of a real-valued function

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad (7.4)$$

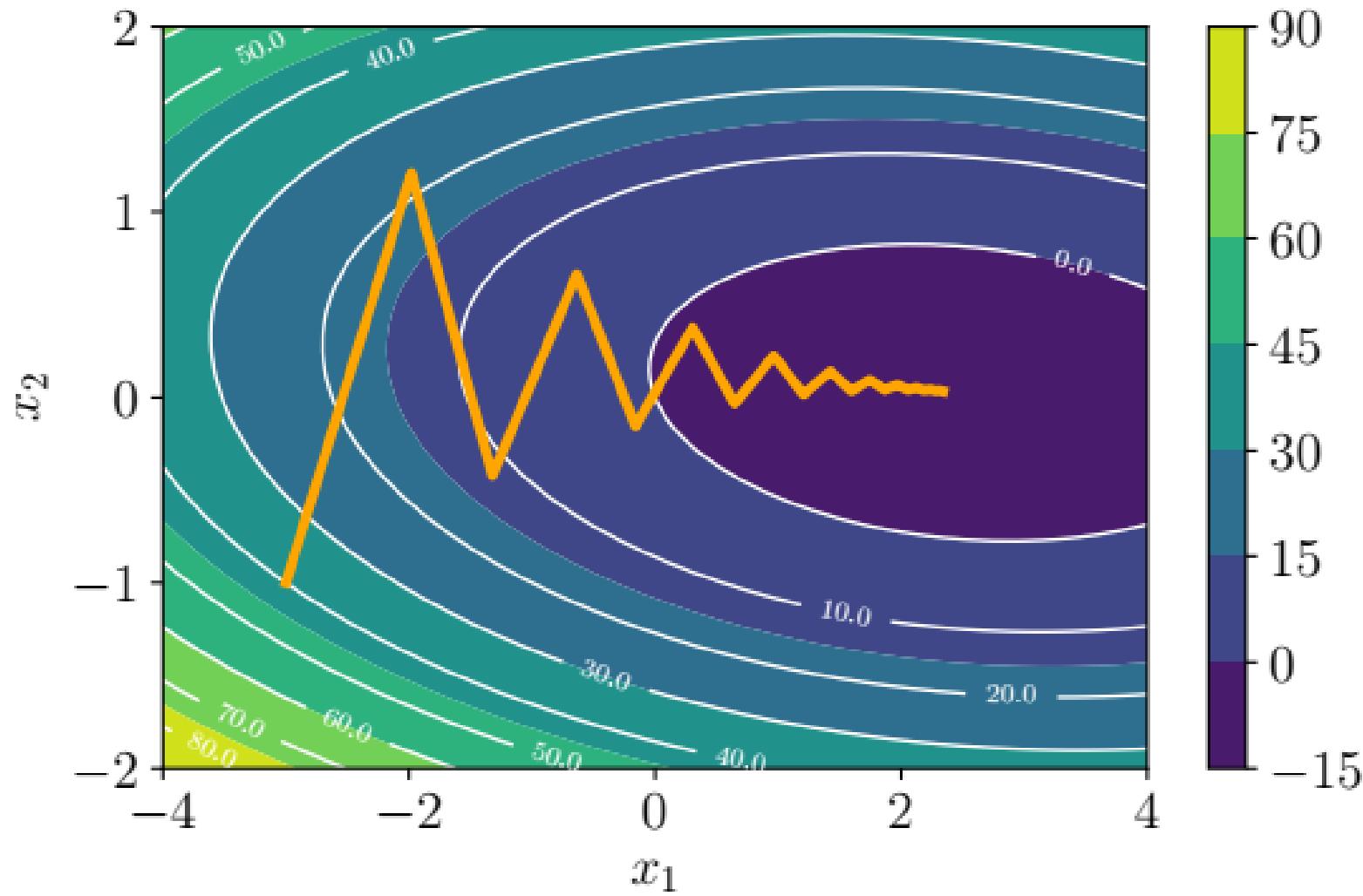
where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is an objective function that captures the machine learning problem at hand. We assume that our function f is differentiable, and we are unable to analytically find a solution in closed form.

Gradient descent is a first-order optimization algorithm. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point.

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i ((\nabla f)(\mathbf{x}_i))^{\top}. \quad (7.6)$$

For suitable step-size γ_i , the sequence $f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq \dots$ converges to a local minimum.

Gradient Descent



Stochastic Gradient Descent

In machine learning, given $n = 1, \dots, N$ data points, we often consider objective functions that are the sum of the losses L_n incurred by each example n . In mathematical notation, we have the form

$$L(\boldsymbol{\theta}) = \sum_{n=1}^N L_n(\boldsymbol{\theta}), \quad (7.13)$$

Standard gradient descent, as introduced previously, is a “batch” optimization method, i.e., optimization is performed using the full training set by updating the vector of parameters according to

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \gamma_i (\nabla L(\boldsymbol{\theta}_i))^T = \boldsymbol{\theta}_i - \gamma_i \sum_{n=1}^N (\nabla L_n(\boldsymbol{\theta}_i))^T \quad (7.15)$$

Consider the term $\sum_{n=1}^N (\nabla L_n(\boldsymbol{\theta}_i))$ in (7.15). We can reduce the amount of computation by taking a sum over a smaller set of L_n . In contrast to batch gradient descent, which uses all L_n for $n = 1, \dots, N$, we randomly choose a subset of L_n for mini-batch gradient descent. In the extreme case, we randomly select only a single L_n to estimate the gradient. The

Probabilities

Sum Rule and Product Rule

The first rule, the *sum rule*, states that

$$p(\mathbf{x}) = \begin{cases} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{y} \text{ is discrete} \\ \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} & \text{if } \mathbf{y} \text{ is continuous} \end{cases}, \quad (6.20)$$

where \mathcal{Y} are the states of the target space of random variable Y . This

The second rule, known as the *product rule*, relates the joint distribution to the conditional distribution via

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}). \quad (6.22)$$

Bayes' Rule

In machine learning and Bayesian statistics, we are often interested in making inferences of unobserved (latent) random variables given that we have observed other random variables. Let us assume we have some prior knowledge $p(x)$ about an unobserved random variable x and some relationship $p(y | x)$ between x and a second random variable y , which we can observe. If we observe y , we can use Bayes' theorem to draw some conclusions about x given the observed values of y . *Bayes' theorem* (also *Bayes' rule* or *Bayes' law*)

$$\underbrace{p(x | y)}_{\text{posterior}} = \frac{\overbrace{p(y | x)}^{\text{likelihood}} \overbrace{p(x)}^{\text{prior}}}{\underbrace{p(y)}_{\text{evidence}}} \quad (6.23)$$

Bayes' Rule

In (6.23), $p(\mathbf{x})$ is the *prior*, which encapsulates our subjective prior knowledge of the unobserved (latent) variable \mathbf{x} before observing any data. We can choose any prior that makes sense to us, but it is critical to ensure that the prior has a nonzero pdf (or pmf) on all plausible \mathbf{x} , even if they are very rare.

The *likelihood* $p(\mathbf{y} | \mathbf{x})$ describes how \mathbf{x} and \mathbf{y} are related, and in the case of discrete probability distributions, it is the probability of the data \mathbf{y} if we were to know the latent variable \mathbf{x} . Note that the likelihood is not a distribution in \mathbf{x} , but only in \mathbf{y} . We call $p(\mathbf{y} | \mathbf{x})$ either the “likelihood of \mathbf{x} (given \mathbf{y})” or the “probability of \mathbf{y} given \mathbf{x} ” but never the likelihood of \mathbf{y} (MacKay, 2003).

The *posterior* $p(\mathbf{x} | \mathbf{y})$ is the quantity of interest in Bayesian statistics because it expresses exactly what we are interested in, i.e., what we know about \mathbf{x} after having observed \mathbf{y} .

Expected Value

Definition 6.3 (Expected Value). The *expected value* of a function $g : \mathbb{R} \rightarrow \mathbb{R}$ of a univariate continuous random variable $X \sim p(x)$ is given by

$$\mathbb{E}_X[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx. \quad (6.28)$$

Correspondingly, the expected value of a function g of a discrete random variable $X \sim p(x)$ is given by

$$\mathbb{E}_X[g(x)] = \sum_{x \in \mathcal{X}} g(x)p(x), \quad (6.29)$$

where \mathcal{X} is the set of possible outcomes (the target space) of the random variable X .

Covariance and Variance

Definition 6.6 (Covariance (Multivariate)). If we consider two multivariate random variables X and Y with states $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^E$ respectively, the *covariance* between X and Y is defined as

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^\top = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}. \quad (6.37)$$

Definition 6.7 (Variance). The *variance* of a random variable X with states $\mathbf{x} \in \mathbb{R}^D$ and a mean vector $\boldsymbol{\mu} \in \mathbb{R}^D$ is defined as

$$\mathbb{V}_X[\mathbf{x}] = \text{Cov}_X[\mathbf{x}, \mathbf{x}] \quad (6.38a)$$

$$= \mathbb{E}_X[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = \mathbb{E}_X[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}_X[\mathbf{x}]\mathbb{E}_X[\mathbf{x}]^\top \quad (6.38b)$$

$$= \begin{bmatrix} \text{Cov}[x_1, x_1] & \text{Cov}[x_1, x_2] & \dots & \text{Cov}[x_1, x_D] \\ \text{Cov}[x_2, x_1] & \text{Cov}[x_2, x_2] & \dots & \text{Cov}[x_2, x_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[x_D, x_1] & \dots & \dots & \text{Cov}[x_D, x_D] \end{bmatrix}. \quad (6.38c)$$

Empirical Mean and Covariance

Definition 6.9 (Empirical Mean and Covariance). The *empirical mean* vector is the arithmetic average of the observations for each variable, and it is defined as

$$\bar{\mathbf{x}} := \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad (6.41)$$

where $\mathbf{x}_n \in \mathbb{R}^D$.

Similar to the empirical mean, the *empirical covariance* matrix is a $D \times D$ matrix

$$\Sigma := \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top. \quad (6.42)$$

Sums of Random Variables

Consider two random variables X, Y with states $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$. Then:

$$\mathbb{E}[\mathbf{x} + \mathbf{y}] = \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{y}] \quad (6.46)$$

$$\mathbb{E}[\mathbf{x} - \mathbf{y}] = \mathbb{E}[\mathbf{x}] - \mathbb{E}[\mathbf{y}] \quad (6.47)$$

$$\mathbb{V}[\mathbf{x} + \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] + \text{Cov}[\mathbf{x}, \mathbf{y}] + \text{Cov}[\mathbf{y}, \mathbf{x}] \quad (6.48)$$

$$\mathbb{V}[\mathbf{x} - \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] - \text{Cov}[\mathbf{x}, \mathbf{y}] - \text{Cov}[\mathbf{y}, \mathbf{x}]. \quad (6.49)$$

Transformations of Random Variables

Mean and (co)variance exhibit some useful properties when it comes to affine transformation of random variables. Consider a random variable X with mean μ and covariance matrix Σ and a (deterministic) affine transformation $\mathbf{y} = A\mathbf{x} + \mathbf{b}$ of \mathbf{x} . Then \mathbf{y} is itself a random variable whose mean vector and covariance matrix are given by

$$\mathbb{E}_Y[\mathbf{y}] = \mathbb{E}_X[A\mathbf{x} + \mathbf{b}] = A\mathbb{E}_X[\mathbf{x}] + \mathbf{b} = A\mu + \mathbf{b}, \quad (6.50)$$

$$\mathbb{V}_Y[\mathbf{y}] = \mathbb{V}_X[A\mathbf{x} + \mathbf{b}] = \mathbb{V}_X[A\mathbf{x}] = A\mathbb{V}_X[\mathbf{x}]A^\top = A\Sigma A^\top, \quad (6.51)$$

respectively. Furthermore,

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}(A\mathbf{x} + \mathbf{b})^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[A\mathbf{x} + \mathbf{b}]^\top \quad (6.52a)$$

$$= \mathbb{E}[\mathbf{x}]\mathbf{b}^\top + \mathbb{E}[\mathbf{x}\mathbf{x}^\top]A^\top - \mu\mathbf{b}^\top - \mu\mu^\top A^\top \quad (6.52b)$$

$$= \mu\mathbf{b}^\top - \mu\mathbf{b}^\top + (\mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mu\mu^\top)A^\top \quad (6.52c)$$

$$\stackrel{(6.38b)}{=} \Sigma A^\top, \quad (6.52d)$$

where $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mu\mu^\top$ is the covariance of X .

Statistical Independence

Definition 6.10 (Independence). Two random variables X, Y are *statistically independent* if and only if

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}). \quad (6.53)$$

Intuitively, two random variables X and Y are independent if the value of \mathbf{y} (once known) does not add any additional information about \mathbf{x} (and vice versa). If X, Y are (statistically) independent, then

- $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{y})$
- $p(\mathbf{x} | \mathbf{y}) = p(\mathbf{x})$
- $\mathbb{V}_{X,Y}[\mathbf{x} + \mathbf{y}] = \mathbb{V}_X[\mathbf{x}] + \mathbb{V}_Y[\mathbf{y}]$
- $\text{Cov}_{X,Y}[\mathbf{x}, \mathbf{y}] = \mathbf{0}$

Definition 6.11 (Conditional Independence). Two random variables X and Y are *conditionally independent* given Z if and only if

$$p(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{y} | \mathbf{z}) \quad \text{for all } \mathbf{z} \in \mathcal{Z}, \quad (6.55)$$

where \mathcal{Z} is the set of states of random variable Z . We write $X \perp\!\!\!\perp Y | Z$ to denote that X is conditionally independent of Y given Z .

Gaussian Distribution

For a univariate random variable, the Gaussian distribution has a density that is given by

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right). \quad (6.62)$$

The *multivariate Gaussian distribution* is fully characterized by a *mean vector* μ and a *covariance matrix* Σ and defined as

$$p(\mathbf{x} | \mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right), \quad (6.63)$$

where $\mathbf{x} \in \mathbb{R}^D$. We write $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mu, \Sigma)$ or $X \sim \mathcal{N}(\mu, \Sigma)$. Fig-

Gaussians are widely used in statistical estimation and machine learning as they have closed-form expressions for marginal and conditional distributions.