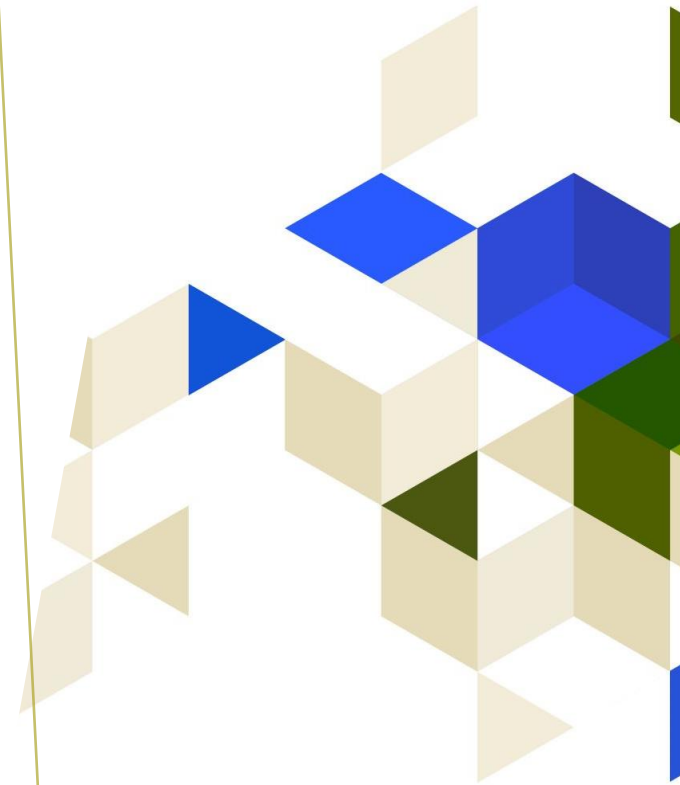


ARTIFICIAL INTELLIGENCE

WOOHWAN JUNG

# Linear Algebra & Probability Theory



# Today's Lecture

- Topics
  - Linear algebra
  - Probability theory
- Not a comprehensive study on linear algebra and probability theory
- Focused on the subset that is most relevant to deep learning

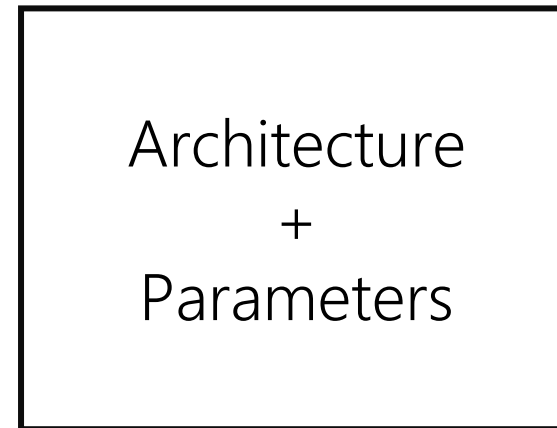
# Learning Process

A lot of vector/matrix operations

Input



Model



Output

$$P(Cat) = 0.9$$

Label

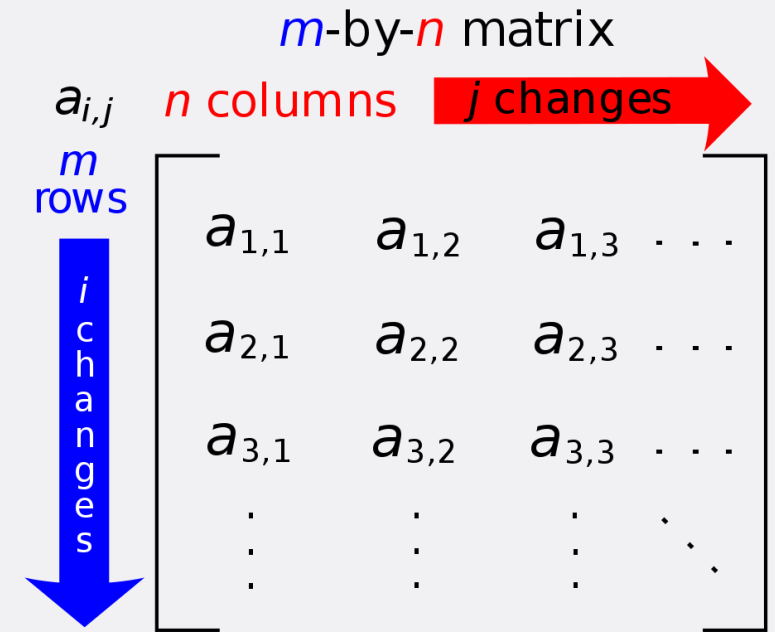
$$P(Cat) = 1$$

Update parameters  
to minimize the loss

Compute loss



# ***LINEAR ALGEBRA***



# Scalars and Vectors

- Scalars
  - A scalar is a single number
  - Usually denoted with italic font:  $a, n, x$
  - Integers, real numbers, rational numbers, etc.
  - Example notation:  $x \in \mathbb{R}, x \in \mathbb{Z}, x \in \mathbb{N}$
- Vectors
  - A vector is a 1-D array of numbers
  - Can be real, integer, etc.
  - Example notation for type and size:
    - $\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{Z}^n$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

# Matrices

- A matrix is a 2-D array of numbers:

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

- Example notation for type and shape:

$$\mathbf{A} \in \mathbb{R}^{m \times n}$$

# Tensors

- A tensor is an (multi-dimensional) array of numbers, that may have
  - Zero dimensions, and be a scalar
  - One dimension, and be a vector
  - Two dimensions, and be a matrix
  - And more dimensions

# Matrix Transpose



$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow \mathbf{A}^T = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

Figure 2.1: The transpose of the matrix can be thought of as a mirror image across the main diagonal.



# Matrix Addition and Subtraction

- Adding or subtracting corresponding elements

- Let  $\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$ ,  $\mathbf{B} = \begin{bmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{bmatrix}$

- $\mathbf{A} + \mathbf{B} = \begin{bmatrix} A_{1,1} + B_{1,1} & A_{1,2} + B_{1,2} \\ A_{2,1} + B_{2,1} & A_{2,2} + B_{2,2} \end{bmatrix}$

- $\mathbf{A} - \mathbf{B} = \begin{bmatrix} A_{1,1} - B_{1,1} & A_{1,2} - B_{1,2} \\ A_{2,1} - B_{2,1} & A_{2,2} - B_{2,2} \end{bmatrix}$

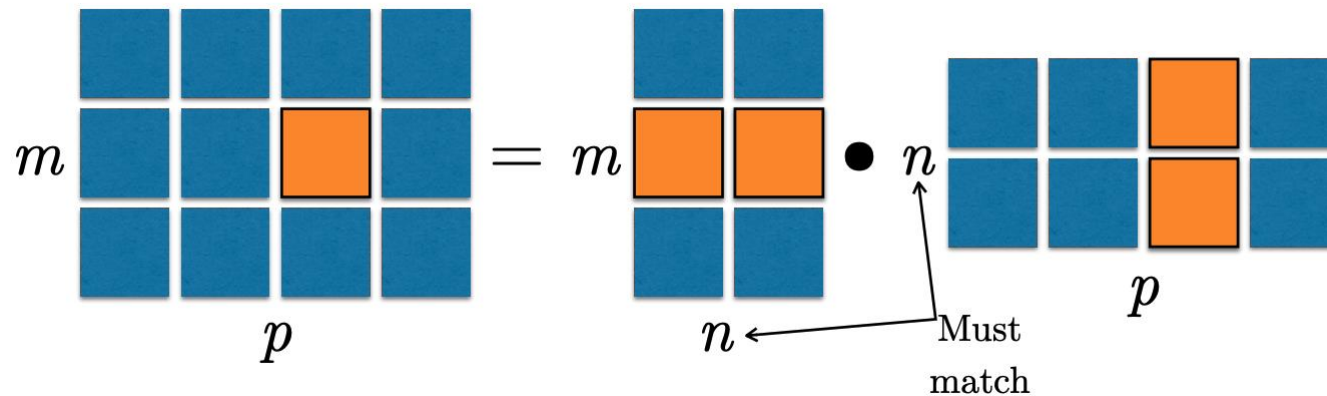
# Matrix Multiplication

- For  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$

$$C = AB.$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}.$$

$$C \in \mathbb{R}^{m \times p}$$



# Shape of the Result of Vector/Matrix Multiplication

- For  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ , and  $C = AB$ 
  - $C \in \mathbb{R}^{m \times p}$
- For  $A \in \mathbb{R}^{m \times n}$ ,  $x \in \mathbb{R}^n$ , and  $y = Ax$ 
  - $y \in \mathbb{R}^m$

# Norms

- Functions that measure how “Large” a vector is
  - Lp norm

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

- Most popular norm: L2 norm ( $p=2$ )
  - L1 norm ( $p=1$ ):  $\|\mathbf{x}\|_1 = \sum_i |x_i|$
  - Max norm ( $p = \infty$ ):  $\|\mathbf{x}\|_\infty = \max_i |x_i|$
- Frobenius norm of a matrix
  - $\|\mathbf{A}\|_F = \sqrt{\sum_i \sum_{ij} A_{ij}^2}$

# Distance Between a Pair of Vectors

- Norm of  $\mathbf{x} - \mathbf{y}$
- Lp distance

$$\|\mathbf{x} - \mathbf{y}\|_p = \left( \sum_i |x_i - y_i|^p \right)^{\frac{1}{p}}$$

# ***PROBABILITY THEORY***

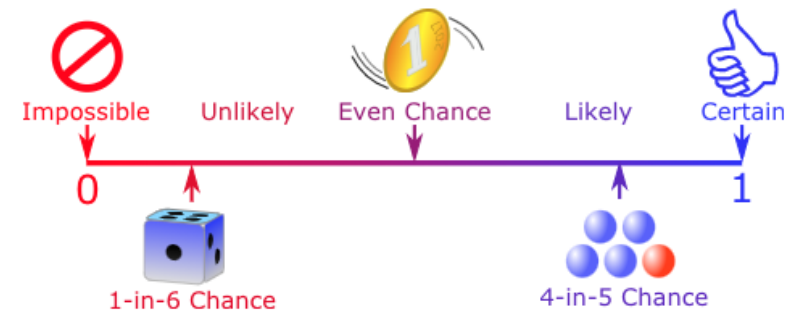


# Random Variable

- A **variable** whose values depend on outcomes of a **random** phenomenon
- Discrete random variables
  - Bernoulli r.v.
  - Categorical r.v.
- Continuous random variables
  - Gaussian (Normal) r.v.
  - Laplace r.v.

# Probability Mass Function (PMF): $P(x)$

- A function that gives the probability that a **discrete random variable** is equal to some value
- For all  $x$ ,  $0 \leq P(X = x) \leq 1$
- $\sum P(X = x) = 1$
- Example: discrete uniform distribution
  - $P(X = x) = \frac{1}{k}$  where  $k$  is the number of possible values

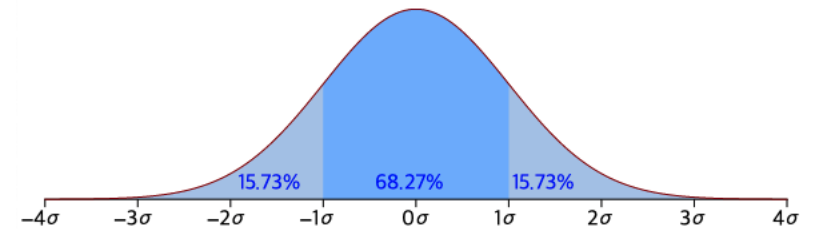


Probability is always between 0 and 1



# Probability Density Function (PDF): $p(x)$

- The PDF of a **continuous random variable** gives the *relative* likelihood of any outcome  $x$
- Properties
  - $p(x) \geq 0$ , (Note: we do not require  $p(x) \leq 1$ )
  - $\int p(x)dx = 1$
- Example: continuous uniform distribution  $u(a,b)$



$$p(x) = \frac{1}{b - a}$$

# Marginal Probability Distribution

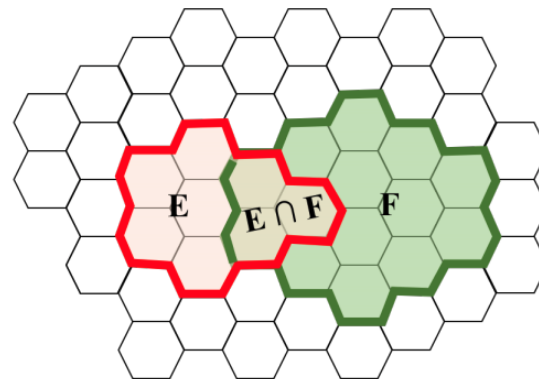
- A probability distribution over the subset of variables
- Sum rule

$$\forall x \in \mathbf{x}, P(\mathbf{x} = x) = \sum_y P(\mathbf{x} = x, y = y).$$

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

# Conditional Probability

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}.$$



$$P(E) = \frac{8}{50} \approx 0.16$$

$$P(E|F) = \frac{3}{14} \approx 0.21$$

# Chain Rule



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

- Example (n=4)

$$\begin{aligned} & P(X^{(1)}, X^{(2)}, X^{(3)}, X^{(4)}) \\ &= P(X^{(1)})P(X^{(2)} \mid X^{(1)})P(X^{(3)} \mid X^{(1)}X^{(2)}) \end{aligned}$$

# Independence

- Two random variables  $x$  and  $y$  are **independent** if and only if

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y).$$

- If  $x$  and  $y$  are **conditionally independent**

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, p(\mathbf{x} = x, \mathbf{y} = y \mid \mathbf{z} = z) = p(\mathbf{x} = x \mid \mathbf{z} = z)p(\mathbf{y} = y \mid \mathbf{z} = z).$$

# Expectation

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

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

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

# Variance & Covariance

- Variance

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

$$= E[f(x)^2] - E[f(x)]^2$$

- Standard deviation: square root of the variance
- Covariance

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

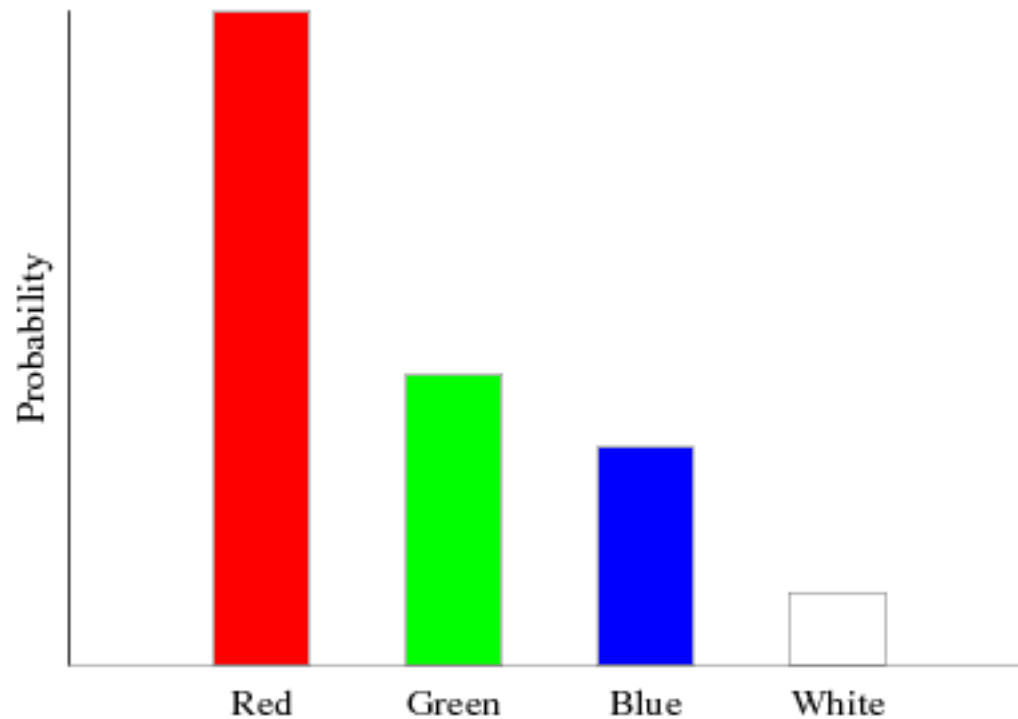
# Bernoulli Distribution

- PMF:  $P(X = x) = \begin{cases} \phi & \text{if } x = 1 \\ 1 - \phi & \text{if } x = 0 \end{cases}$
- Expectation:  $E[X] = \phi$
- Variance:  $Var[X] = \phi(1 - \phi)$



# Categorical Distribution

- A.k.a. multinoulli distribution



$$P(X = k) = p_k$$

# Gaussian Distribution

- A.k.a Normal distribution
- Parameterized by variance  $\sigma^2$ :

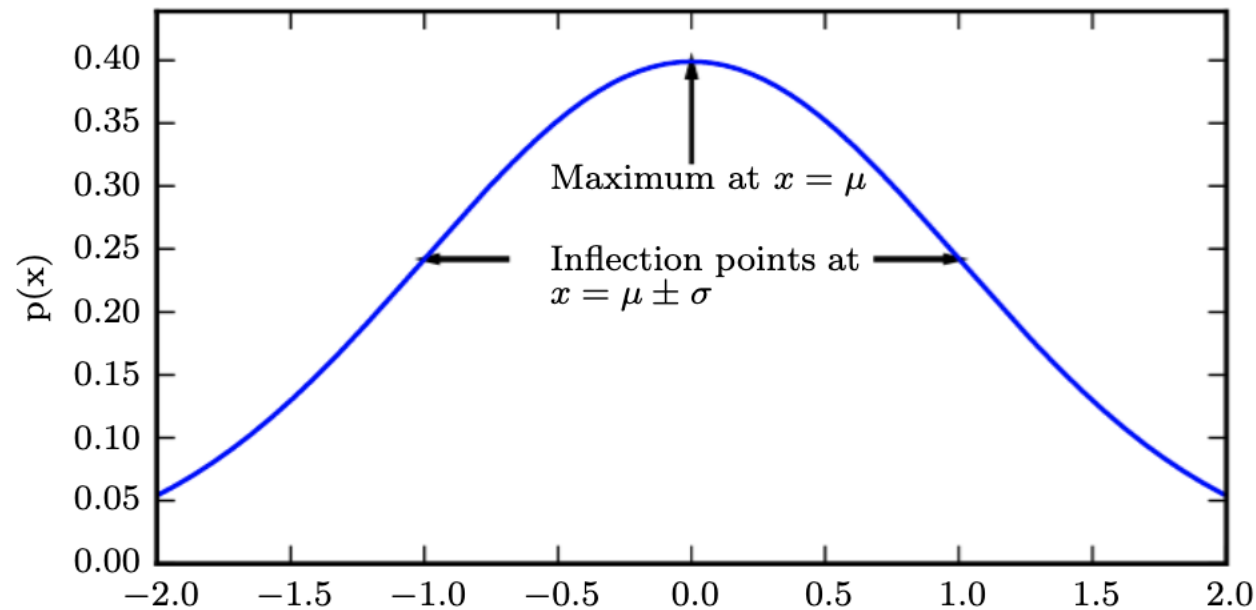
$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right).$$

- Parameterized by precision  $\beta = \frac{1}{\sigma^2}$

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

# Standard Normal Distribution

- Normal (gaussian) distribution  $\mathcal{N}(x; \mu, \sigma^2)$  with  $\mu = 0$  and  $\sigma = 1$

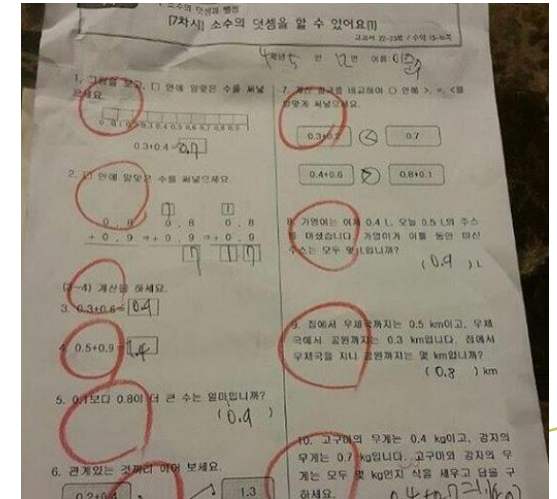


# Gaussian is Indeed Normal!

- Gaussian distributions are sensible choice for many application
  - Especially in the absence of prior knowledge

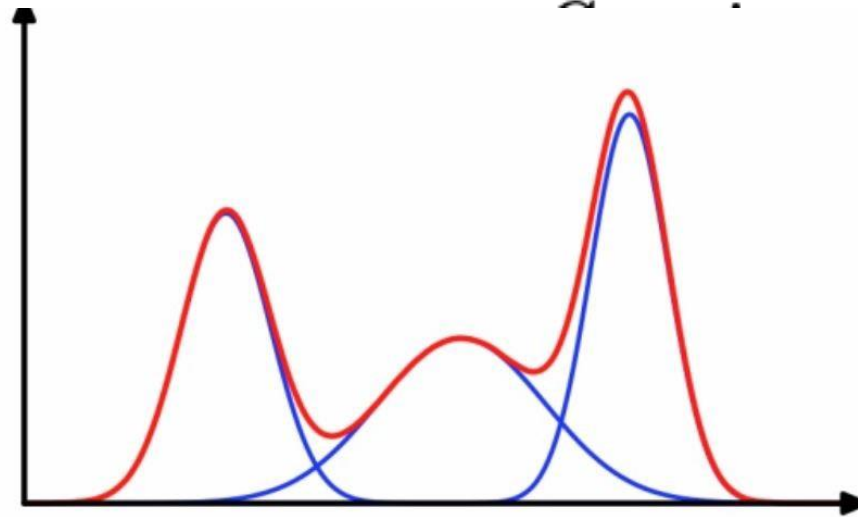
# Central limit theorem

"The sum of many independent random variables is approximately **normally** distributed"



# Mixture of Distributions

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} \mid c = i)$$



Gaussian mixture

# Bayes' Rule

$$\text{Posterior } P(x | y) = \frac{\text{Prior } P(x) \text{ Likelihood } P(y | x)}{P(y)}.$$

# Learning Process

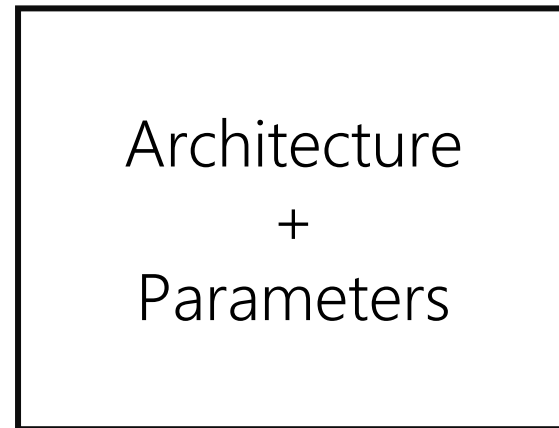
A lot of vector/matrix operations

Input

Model

Output

Label



$$P(Cat) = 0.9$$

$$P(Cat) = 1$$

Update parameters  
to minimize the loss

Compute loss

Q1: How to output probability distributions (from the results of vector/matrix operations)?

Q2: How to measure the distance between two probability distributions?

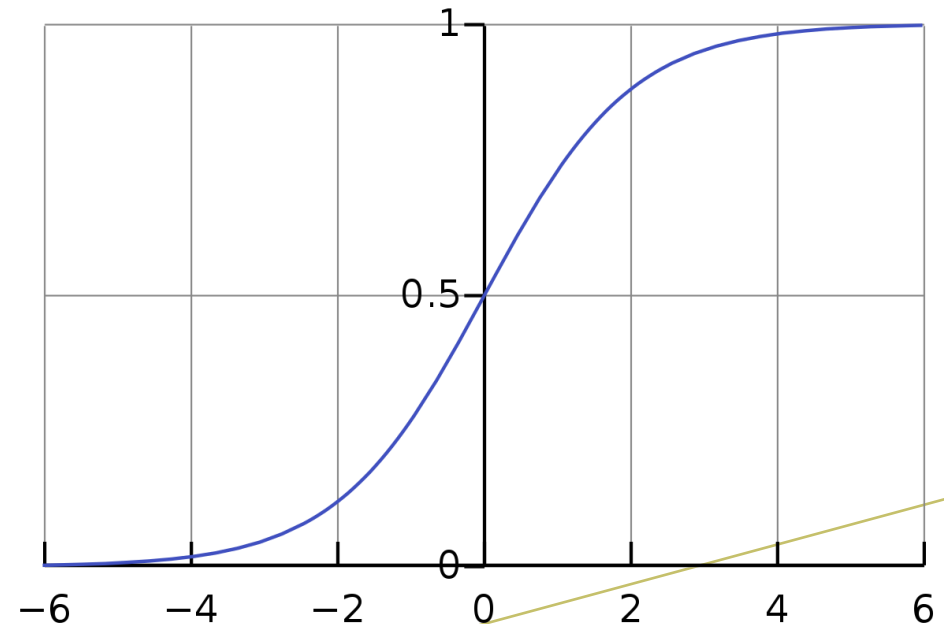
# Q1: How to output probability distributions?

- Range of vector/matrix operations:  $(-\infty, +\infty)$
- Range of  $P(X)$  in Bernoulli distribution  $[0,1]$
- Logistic sigmoid

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

It's not standard deviation

- Domain:  $-\infty \sim +\infty$
- Range:  $0 \sim 1$





# Q1: How to output probability distributions?

- Softplus function

$$\zeta(x) = \log(1 + \exp(x)).$$

- Softened version of

$$x^+ = \max(0, x).$$

- Can be used to represent standard deviations

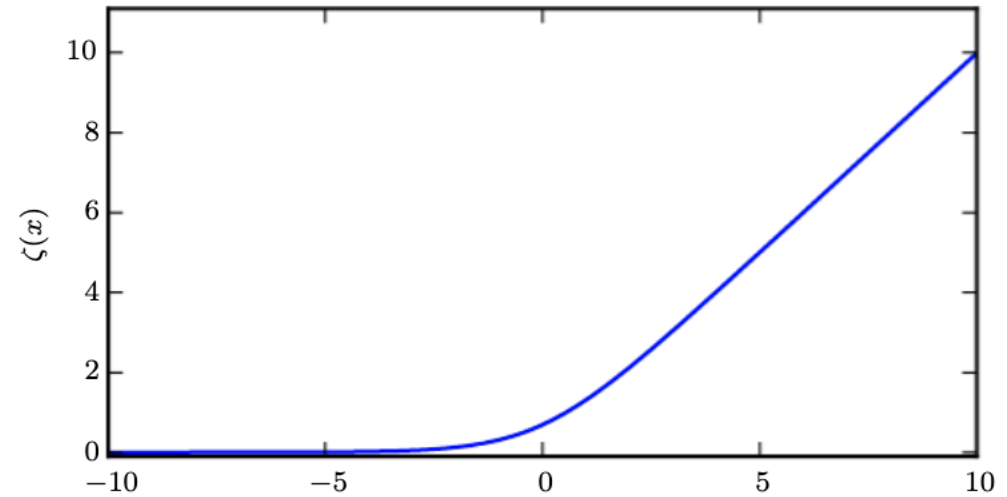


Figure 3.4: The softplus function.

## **Q2: How to Measure the Distance Between Two Probability Distributions?**

A: Information Theory

***INFORMATION  
THEORY  
(OPTIONAL)***

# Information

- Quantity of information

**1000 bits**

00000000...0000000000

Same quantity?



**1000 bits**

0010001...111001001

Same quantity?



0 \* 1000

0\*2, 1\*1, 0\*3...1\*3, 0  
\*2, 1\*1, 0\*2, 1\*1

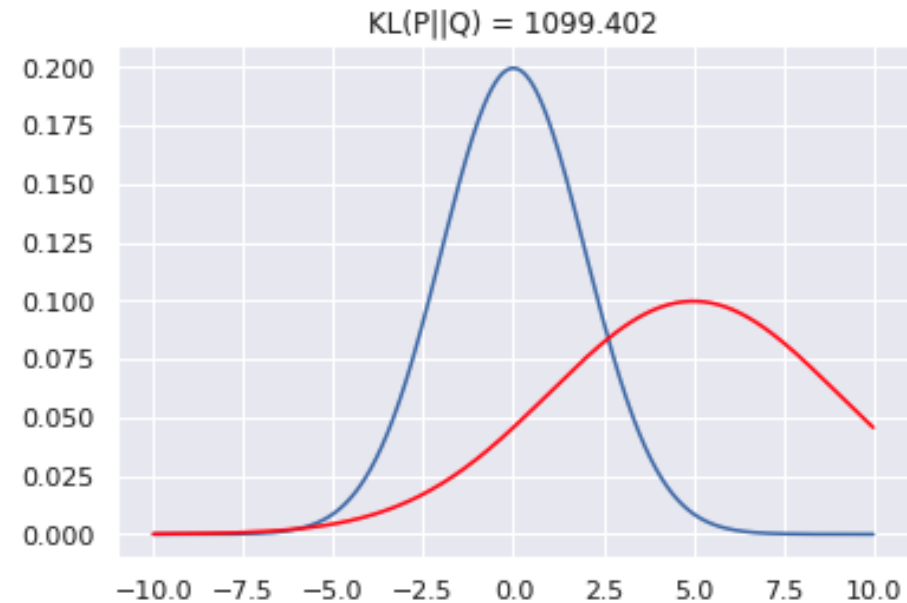
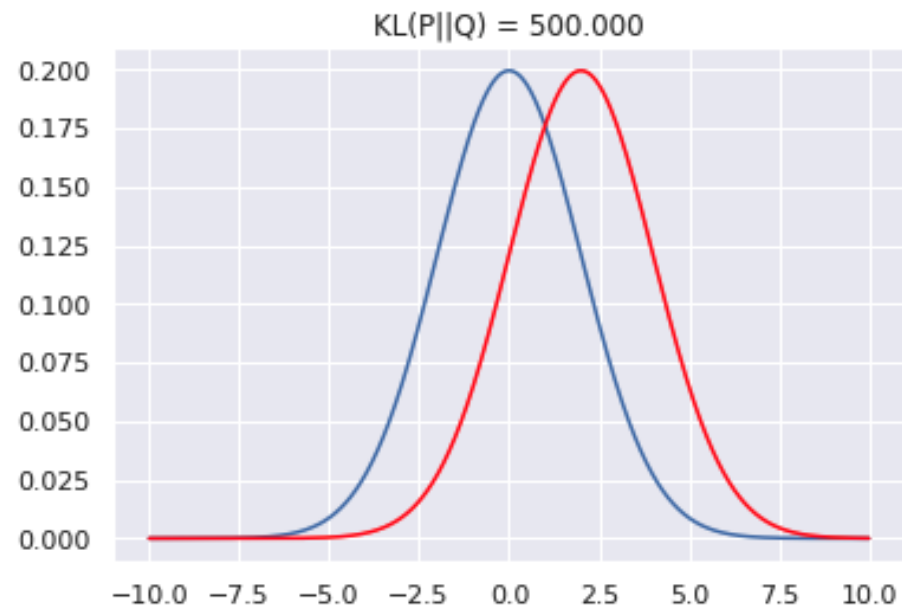
# (Self) Information $I(x)$

- Roughly speaking, the minimum number of bits to encode a signal  $x$
- Definition
  - $I(x) = -\log P(x)$
- Intuition
  - If a pattern is frequent, it can be simply and efficiently encoded/compressed
  - Example: 00000000...0000000000

# Information Theory

- Self information of an event  $x = x$ 
  - $I(x) = -\log P(x)$
- Entropy
  - $H(X) = E_{X \sim P}[I(x)] = -E_{X \sim P}[\log P(x)]$
  - Computation:  $H(X) = -\sum_x P(x) \log P(x)$
- Kullback-Leibler(KL) divergence
  - $D_{KL}(P||Q) = E_{X \sim P} \left[ \log \frac{P(x)}{Q(x)} \right]$ 
    - $P, Q$ : probability distributions
  - Computation:  $D_{KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$

# Kullback-Leibler(KL) Divergence



# Appendix: Why L2 Norm?

- Likelihood:  $P(\text{Observation}|\text{Model})$
- Log-likelihood:  $\log P(\text{Observation}|\text{Model})$

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

- $\log P(X) = -\frac{1}{2\sigma^2}(x - \mu)^2 + C$

Minimizing L2  
distance

$\approx$

Maximizing  $\log P(X)$   
where  $X \sim N$