

SVD & Information Theory

Fartash Faghri

University of Toronto

CSC2515, Fall 2019

Eigen-values & Eigen-vectors

Vectors v and scalars λ that s.t. $\mathbf{A}v = \lambda v$

Col space, $\mathbf{A}v = \lambda v$ Row space, $u^\top \mathbf{A} = \lambda u^\top$

What are the e-values and e-vectors of A, B:

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}$$

Singular Value Decomposition (SVD)

Any real matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ can be decomposed as

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^{\top}$$

where $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{m \times m}$ are orthonormal

and $\mathbf{S} \in \mathbb{R}^{n \times m}$ is diagonal.

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

$$\mathbf{U} \in \mathbb{R}^{n \times n} \quad \text{e-vectors of} \quad \mathbf{A}\mathbf{A}^\top = \mathbf{U}\mathbf{S}^2\mathbf{U}^\top \\ = \mathbf{U}\mathbf{S}\mathbf{V}^\top(\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top = \mathbf{U}\mathbf{S}\mathbf{V}^\top\mathbf{V}\mathbf{S}\mathbf{U}^\top$$

$$\mathbf{V} \in \mathbb{R}^{m \times m} \quad \text{e-vectors of} \quad \mathbf{A}^\top\mathbf{A} = \mathbf{V}\mathbf{S}^2\mathbf{V}^\top$$

$$\mathbf{S} \in \mathbb{R}^{n \times m} \quad \text{Singular-values (non-negative)}$$

*Eigen-value decomposition and SVD are not the same.

*Even if eigen-value decomposition is defined, eigen-values and singular-values are not generally the same.

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

What is the SVD of A, B?

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}$$

Positive Semi-Definite (PSD)

$$Av = \lambda v$$

Definition: $\forall x \in \mathbb{R}^n, \quad xAx^\top \geq 0$

Also means, all e-values are non-negative: $\lambda \geq 0$

E.g. in the normal distribution, the probability is non-negative:

(1D)

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

(nD)

$$\frac{1}{(2\pi)^{\frac{k}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}$$

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

Are these PSD?

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}$$

Matrix Inverse

$$\mathbf{A}\mathbf{A}^{-1} = \mathbb{I}_{n \times n}$$

If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is full rank, $\mathbf{A}^{-1} \in \mathbb{R}^{n \times n}$

If not full rank, inverse doesn't exist (pseudoinverse)

If not square, inverse is not defined:

$$\mathbf{A}^{-1}\mathbf{A} = \mathbb{I}_{m \times m} \quad \text{or} \quad \mathbf{A}\mathbf{A}^{-1} = \mathbb{I}_{n \times n} \quad ?$$

Matrix inverse using SVD $\mathbf{A}\mathbf{A}^{-1} = \mathbb{I}_{n \times n}$ $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$

If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is full rank $\mathbf{A}^{-1} = \mathbf{U}\mathbf{S}^{-1}\mathbf{U}^{-1}$

If not full rank, $\mathbf{A}^+ = \mathbf{V}\mathbf{S}^+\mathbf{U}^\top$

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

Moore-Penrose inverse is the most common pseudoinverse.

Common in practice: $\mathbf{A}^+ = \mathbf{V}(\mathbf{S} + \lambda\mathbb{I})^{-1}\mathbf{U}^\top$ (not MP-inv)

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

What is the inverse?

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}$$

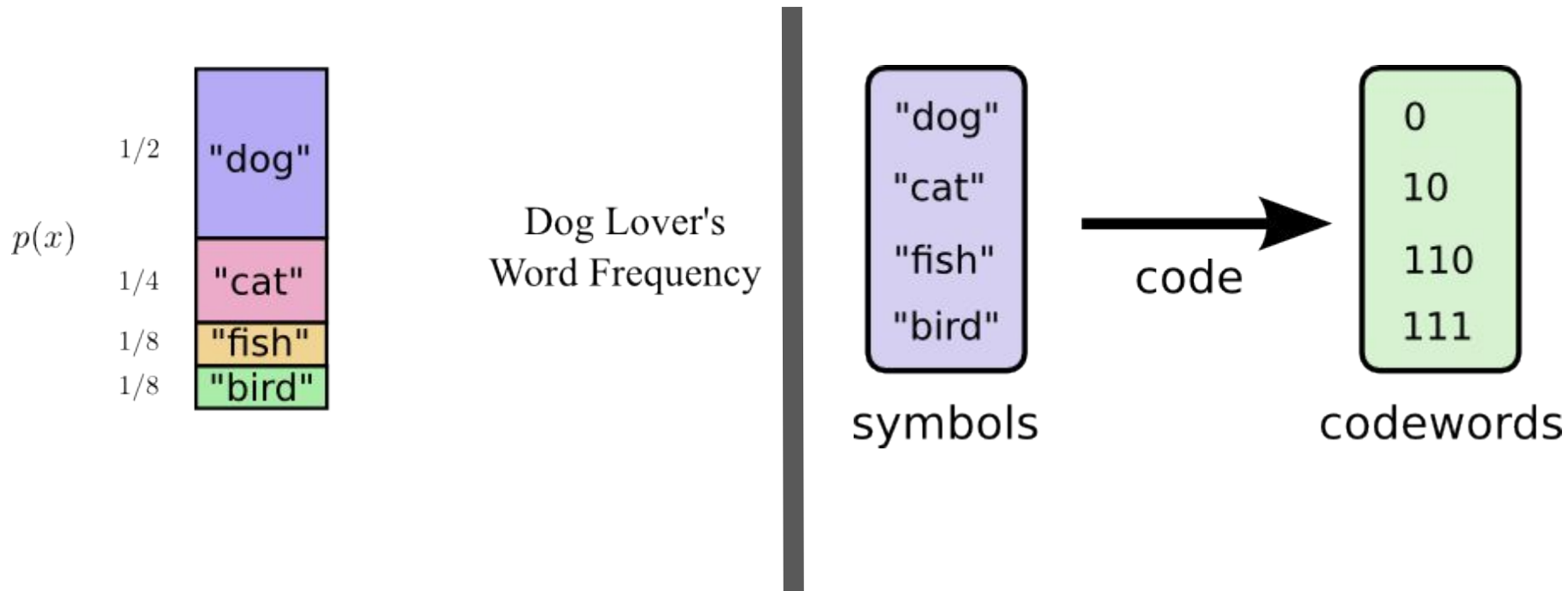
Information theory

You will see it in

- Entropy and Information gain in Decision Trees,
- KL divergence to measure distances between probability distributions,
- cross-entropy loss that is widespread in training classifiers.

Coding

What is the minimum length of code for communicating the messages of a dog lover?



The space of codewords

2 codes of length 1 (0, 1)

4 codes of length 2 (00, 01, 10, 01)

8 codes of length 3 (000, 001, ...)

.....

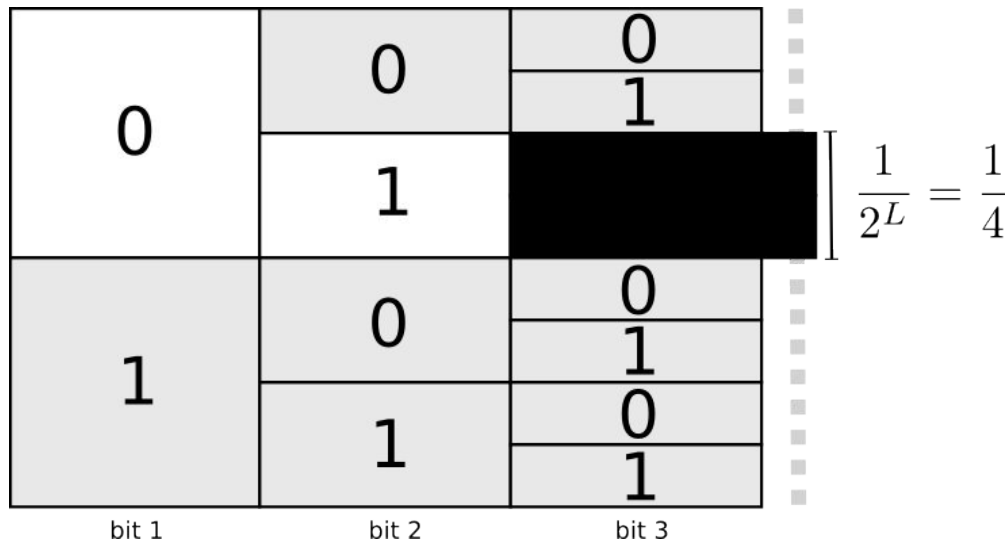
0	0	0
		1
1	1	0
		1
1	0	0
		1
1	1	0
		1

bit 1 bit 2 bit 3



Cost of 01 is the cost of all codes that cannot be used (black)

Cost of length 0 is 1



Greedy decoding:

0110111

Using both codes at the same time is ambiguous:

01 (dog) and 011 (cat)

Optimal Cost

The optimal cost for an event that happens with probability $p(x)$ is $p(x)$ of our total budget.

Pay more for frequently events and less for rare events.

Entropy

Cost of a message of length $L : \frac{1}{2^L}$

Invert to get the length of a message that costs $C : \log_2\left(\frac{1}{C}\right)$

Since we spend $p(x)$ on the codeword for x , it has length $\log_2\left(\frac{1}{p(x)}\right)$

Entropy of a distribution: the average length of the best possible code

$$H(p) = \sum_x p(x) \log_2\left(\frac{1}{p(x)}\right)$$

Entropy is measured in bits (log base-2) or nats (log base-e)

Entropy

$$H(p) = \sum_x p(x) \log_2\left(\frac{1}{p(x)}\right)$$

What is the entropy of each distribution?

- Bernoulli(0)
- Bernoulli(1)
- Bernoulli(0.5)

Information Gain (used in Decision Trees)

Conditional Entropy:
$$H(X|Y) = \sum_{x,y} p(x,y) \log \left(\frac{1}{p(x|y)} \right)$$

How much information I gain by observing X after I had observed Y:

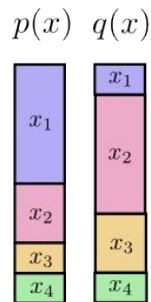
Information Gain:
$$IG(X, Y) = H(X) - H(X|Y)$$

Cross-entropy (used in Classification)

The average length of communicating an event from one distribution with the optimal code for another distribution

Entropy:
$$H(p) = \sum_x p(x) \log \left(\frac{1}{p(x)} \right)$$

Cross-Entropy:
$$H_p(q) = \sum_x q(x) \log \left(\frac{1}{p(x)} \right)$$



Cross-Entropy: $H_p(q)$

Average Length
of message from $q(x)$
using code for $p(x)$.

Classification

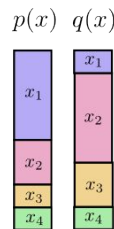
Input data, ground-truth target, (\mathbf{x}_i, t_i) (a single data point i)

Prediction, $y_i \in \{1, \dots, C\}$ (categorical random variable)

Probabilistic classifier, $p(y_i | \mathbf{x}_i; \boldsymbol{\theta})$

Ground-truth distribution, $q(y_i | \mathbf{x}_i) = \begin{cases} 1 & y_i = t_i \\ 0 & y_i \neq t_i \end{cases}$

$$H_p(q) = \sum_x q(x) \log \left(\frac{1}{p(x)} \right)$$



Cross-Entropy: $H_p(q)$

Average Length
of message from $q(x)$
using code for $p(x)$.

Cross-entropy: average length of the ground truth ground-truth using the optimal

code for p :
$$H_p(q) = -\log(p(t_i | \mathbf{x}_i; \boldsymbol{\theta}))$$

Kullback–Leibler (KL) divergence

Distance between two probability distributions

$$D_q(p) = H_q(p) - H(p)$$

$$D_q(p) = \sum_x p(x) \log \left(\frac{p(x)}{q(x)} \right)$$

How much longer our messages are (from p) because we used a code optimized for a different distribution (q). If the distributions are the same, this difference will be zero.

KL Divergence in Classification

Probabilistic classifier, $p(y_i | \mathbf{x}_i; \boldsymbol{\theta})$

Ground-truth distribution, $q(y_i | \mathbf{x}_i) = \begin{cases} 1 & y_i = t_i \\ 0 & y_i \neq t_i \end{cases}$

Cross-entropy, $H_p(q) = \sum_x q(x) \log \left(\frac{1}{p(x)} \right)$

KL Divergence, $D_q(p) = H_q(p) - H(p)$

What is the Entropy of q?