



THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

CS 4364/6364

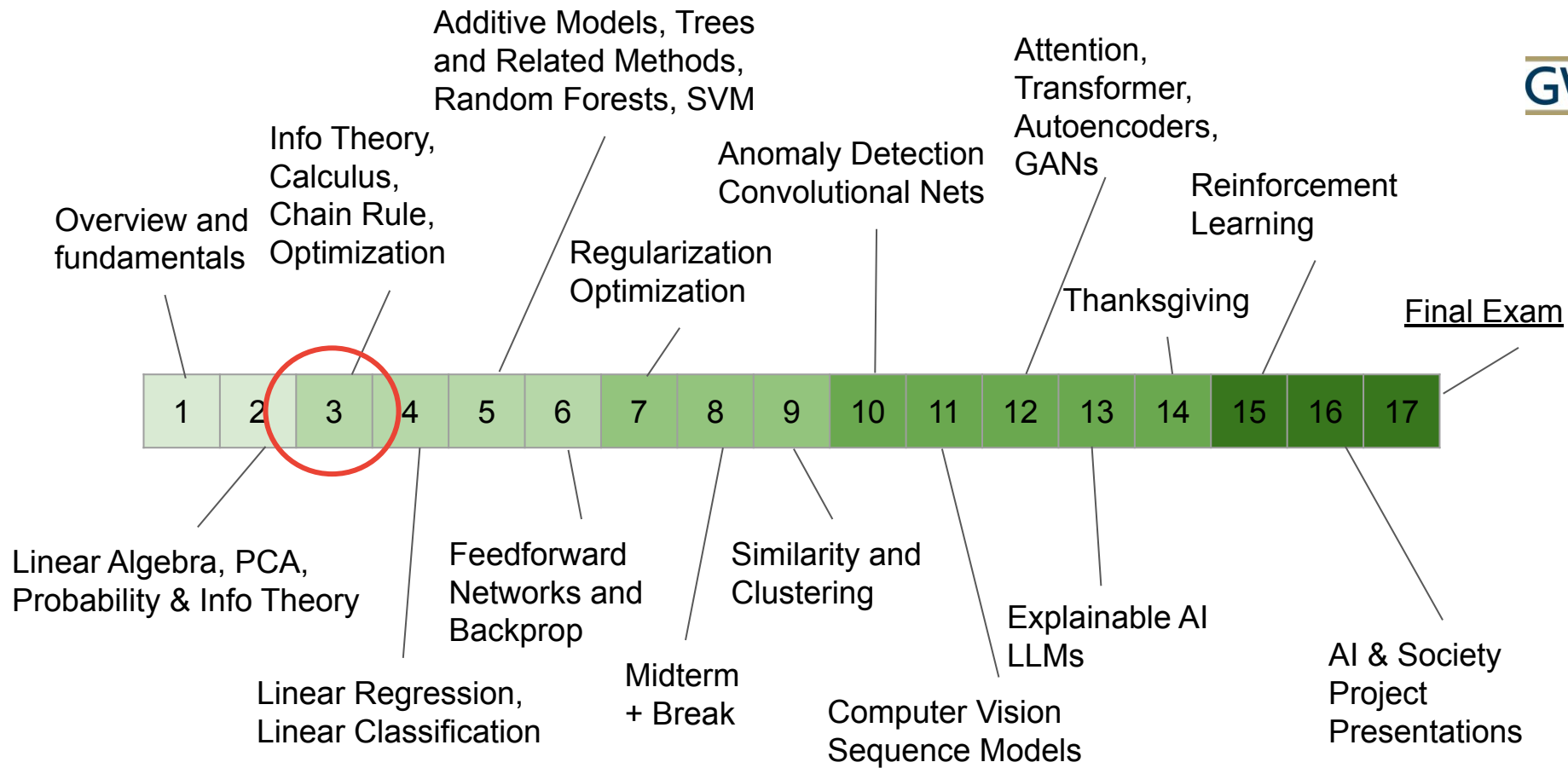
Machine Learning

Fall Semester 9/5/2023

Lecture 4.

From Coin Flips to KL Divergence

John Sipple
jsipple@gwu.edu



Transforming probability distributions

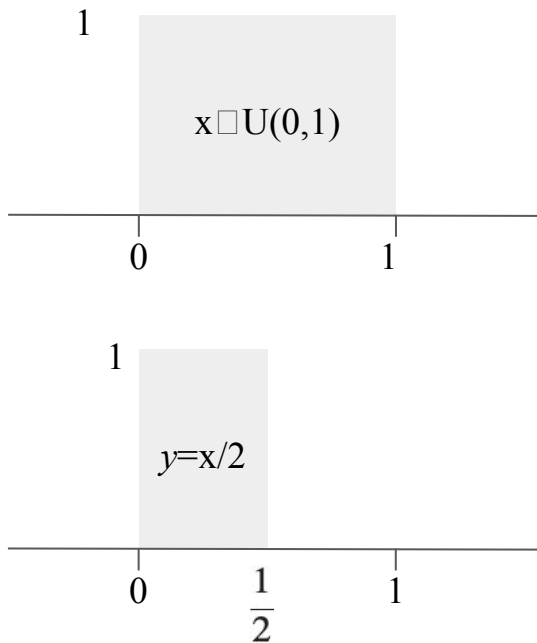
Suppose we have: $y = g(x) = \frac{x}{2}$ and $x \sim U(0, 1)$

And we apply the formula $p_y(y) = p_x(g^{-1}(x))$

Then we violate the rules of probability:

$$\int p_y(y) dy = \frac{1}{2}$$

We need to be cautious about applying functions to probability distributions!



Not a probability distribution!

Change of Variables on PDs requires rescaling

$$y = g(x)$$

$$p_y(y) \neq p_x(g^{-1}(y))$$

We need to preserve the property:

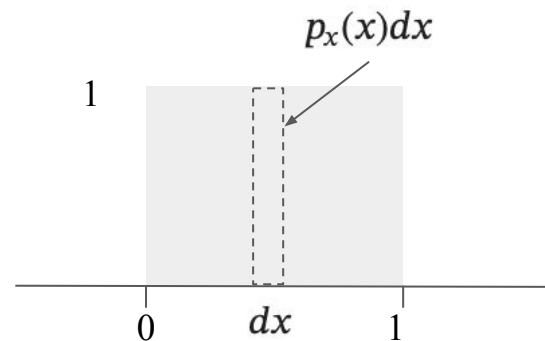
$$|p_y(g(x))dy| = |p_x(x)dx|$$

Can be rewritten as:

$$p_x(x) = p_y(g(x)) \left| \frac{\partial g(x)}{\partial x} \right|$$

And in higher dimension:

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left(\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$



Binomial Random Variable

Each event takes on one of two possible outcomes.

$$E = \{0, 1\}$$

for example, 0 = Tails and 1 = Heads

Specified with a parameter p , where:

$$p = P(E = 1)$$

and

$$1 - p = P(E = 0)$$



i.i.d.

Assume the events are **i.i.d.**:

- **independent:**

$$P(E_1 = 1, E_2 = 1) = P(E_1 = 1)P(E_2 = 1)$$

and

- **identically distributed:** p is constant

$$p_{E_1} = p_{E_2} = \dots$$

Bernoulli Distribution

Describes the probability of seeing s successes in N trials, given a parameter p .

For example, what's the probability of seeing 2 heads in 3 trials with a fair coin, $p = 0.5$?

$$b(s, N; p) =$$

Bernoulli Distribution

Let's enumerate all possible, and exclusive outcomes:

Outcome a. 1, 1, 0: $P(E = 1) \times P(E = 1) \times P(E = 0) = p^2(1 - p)$

Outcome b. 0, 1, 1: $P(E = 0) \times P(E = 1) \times P(E = 1) = p^2(1 - p)$

Outcome c. 1, 0, 1: $P(E = 1) \times P(E = 0) \times P(E = 1) = p^2(1 - p)$

Since the outcomes are exclusive (i.e., a or b or c) could occur, we can add them:

$$P(a \cup b \cup c) = P(a) + P(b) + P(c)$$

Then

$$b(2, 3; 0.5) = 3 \times p^2(1 - p) = 3 \times 0.5^3 = 0.375$$

Bernoulli Distribution (general form)

In general, we'll just count all possible permutations of s successes in N **exclusive** events.

$\binom{N}{s}$: The number of possible permutations of having s successes in N trials.

$$\binom{N}{s} = \frac{N!}{(N-s)!s!}$$

Back to the example:

$$\binom{3}{2} = \frac{3!}{(3-2)!2!} = \frac{1 \times 2 \times 3}{1 \times 1 \times 2} = 3$$

In general, the Binomial distribution is:

$$b(s, N; p) = \binom{N}{s} p^s (1-p)^{N-s}$$

Bernoulli Distribution with multiple levels

Given p that you'll be exposed to the Coronavirus and catch Covid-19 on a Metro ride.

One might ask, what's the probability that you'll be exposed once **or more** in $N = 6$ trips?

$$B(s > 1, N = 6; p) = b(s = 1, N; p) + b(s = 2, N; p) + \cdots = \sum_{s=1}^6 b(s, N; p)$$

Alternatively:

$$B(s > 1, N = 6; p) = 1 - b(s = 0, N; p) = 1 - (1 - p)^6$$

In other words, we can restate the question as, 1 minus the probability of zero exposures in 6 trips.

Information Theory

The study of the transmission, processing, extraction, and utilization of information.

Formalized in 1948 by Claude Shannon: [A Mathematical Theory of Communication](#)

- Reconstruct messages accurately despite a noisy channel

Foundational to data compression and digital communication



Information Theory

Basic Idea: Learning that an unlikely event occurred is more informative than learning that a likely event occurred.

Information of a single event: $I(x) = -\log P(\mathbf{x} = x)$

Entropy of a distribution:

$$H(\mathbf{x}) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\log P(x)]$$

$$H(\mathbf{x}) = -\sum_x P(\mathbf{x} = x) \log P(\mathbf{x} = x)$$

Entropy

The universal measure of **randomness**.

For a discrete random variable x , entropy is defined as:

$$\begin{aligned} H(x) &= - \sum_x P(x = x) \log P(x = x) \\ &= \mathbb{E} [-\log P(x = x)] \end{aligned}$$

Note on log bases: Where the log base can be e , 2, 10 etc., but if it's $|x|$, then maximum entropy, $\max H(x) = 1$. (The shape of H is the same for all log bases.)

Entropy

For two outcomes (i.e., coin flips):

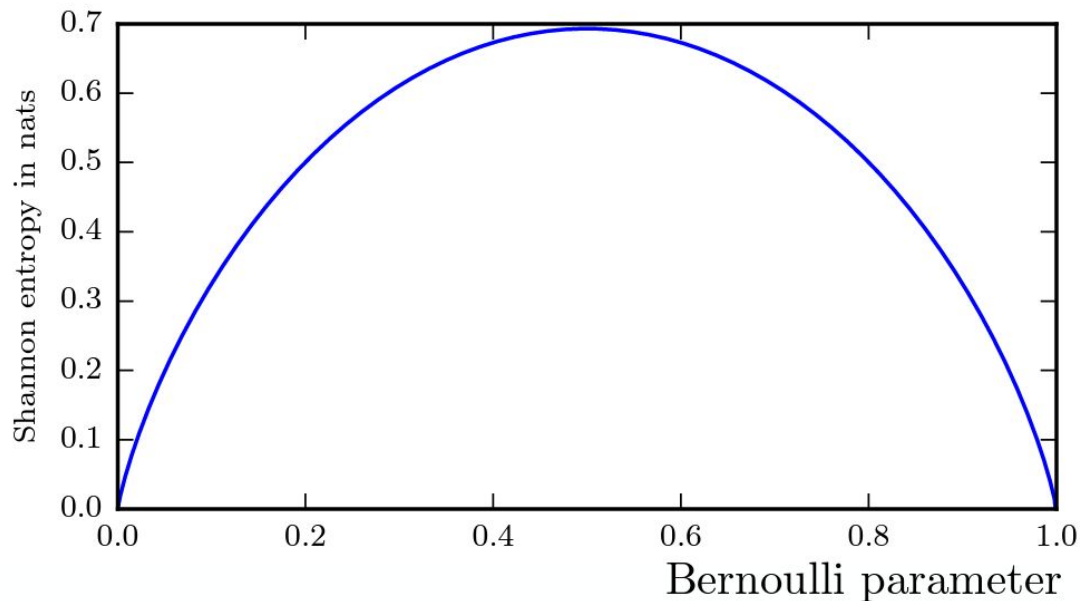
$$H(p) = -p \log p - (1 - p) \log(1 - p)$$

For example, $p = 0.4$:

$$H(p = 0.4) = -0.4 \log 0.4 - 0.6 \log 0.6 = 0.673$$

How about multiple outcomes (e.g, die)?

Entropy of a Coin Flip



$$H(p) = -p \log p - (1 - p) \log(1 - p)$$

Cross-Entropy of two probability distributions

In Machine Learning we are often interested in approximating a true, but unknown RV (p) with an estimator function (q)

Cross-entropy can be a good measure for how well q approximates p .

$$H(p, q) = - \sum_i p \log q$$

For the coin-flip case:

$$H(p, q) = -p \log q - (1 - p) \log(1 - q)$$

Cross-Entropy of two probability distributions

For example:

- **True probability** $p = 0.4$
- **Estimated probability** $q = 0.5$

$$H(p = 0.4, q = 0.5) = -0.4 \log 0.5 - 0.6 \log 0.5 = 0.693$$

Which is larger than $H(p = 0.4) = 0.673$ and is minimized when $q = p$.

Kullback-Leibler (KL) Divergence

A measure of how one probability distribution q is different from another probability distribution p .

$$\begin{aligned} D_{KL}(p||q) &= H(p, q) - H(p) \\ &= - \sum_i p \log q - \left(- \sum_i p \log p \right) \\ &= - \sum_i p \log \frac{q}{p} \end{aligned}$$

Note:

$$D_{KL}(p||p) = 0$$

$$D_{KL}(p||q) \neq D_{KL}(q||p)$$

Kullback-Leibler (KL) Divergence

Back to the example, $p = 0.4$ and $q = 0.5$:

$$D_{KL} (p = 0.4 || q = 0.5) = 0.693 - 0.673 = 0.02$$

Kullback-Leibler (KL) Divergence

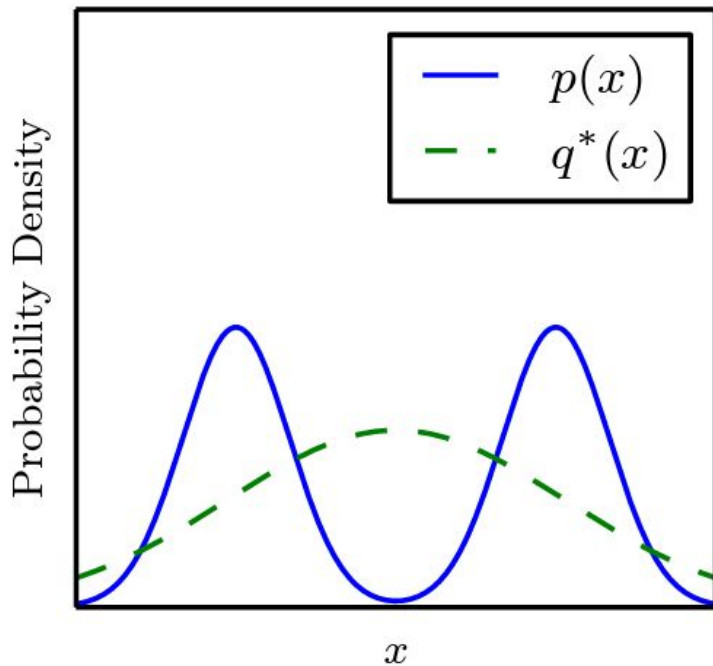
Measures the difference between two distributions from the “perspective” of one distribution

$$D_{KL}(P||Q) = \mathbb{E}_{\mathbf{x} \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{\mathbf{x} \sim P} [\log P(x) - \log Q(x)]$$

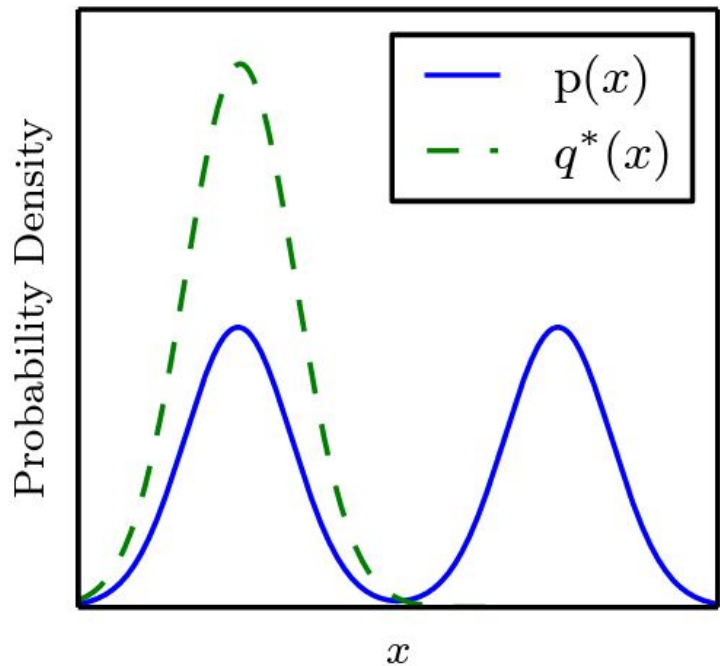
$$D_{KL}(P||Q) = \sum_x P(\mathbf{x} = x) \log \frac{P(\mathbf{x} = x)}{Q(\mathbf{x} = x)}$$

KL Divergence is Asymmetric

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p \| q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q \| p)$$



Binary Cross-Entropy Loss

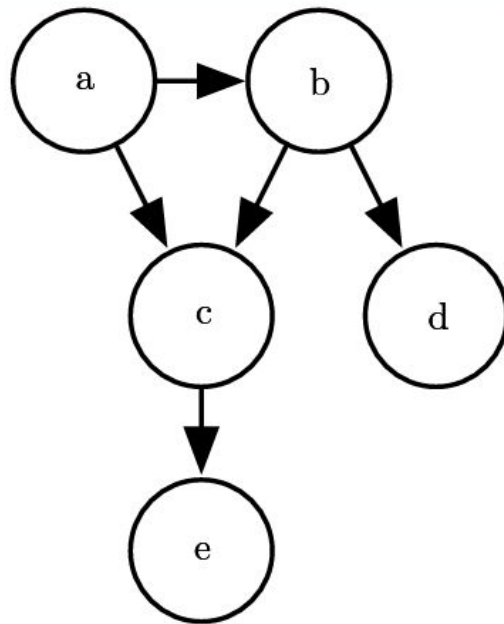
Suppose you want to train a binary classifier and you need a loss function that you want to minimize.

- **Test Data:** $p = [1, 1, 0, 1, 0, 1, \dots]$
- **Classifier predictions:** $q = [0.8, 0.9, 0.3, 0.7, 0.1, 0.9, \dots]$

$$\mathcal{L}(p, q) = \frac{1}{N} \sum_i^N -p_i \log q_i - (1 - p_i) \log(1 - q_i)$$

Directed Structured Model

Figure 3.7



$$p(a, b, c, d, e) = p(a)p(b \mid a)p(c \mid a, b)p(d \mid b)p(e \mid c).$$

Undirected Structured Model

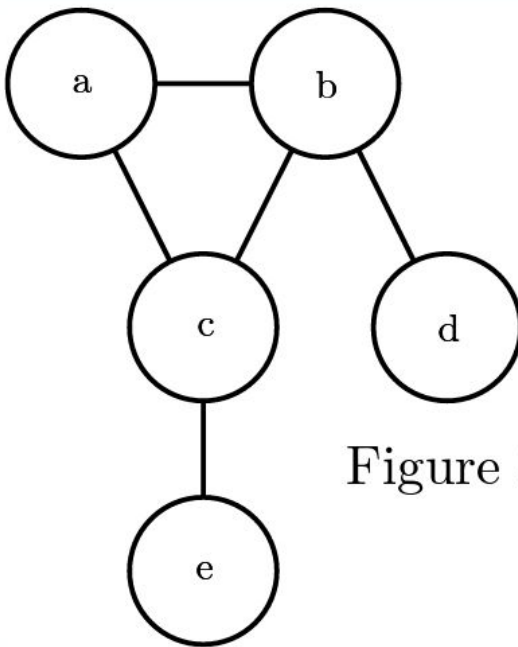


Figure 3.8

$$p(a, b, c, d, e) = \frac{1}{Z} \phi^{(1)}(a, b, c) \phi^{(2)}(b, d) \phi^{(3)}(c, e).$$