National Ph.D. Program in *Artificial Intelligence for Society* **Statistics for Machine Learning** Lesson 04 - Functions of random variables. Distances between distributions. Simulation.

Andrea Pugnana, Salvatore Ruggieri

Department of Computer Science University of Pisa, Italy andrea.pugnana@di.unipi.it salvatore.ruggieri@unipi.it

Functions of two or more random variables: expectation

- $V = \pi H R^2$ be the volume of a vase of height H and radius R
- $g(H, R) = \pi H R^2$ is a random variable (function of random variables)
- $P_V(V=3) = P_{HR}(\pi HR^2 = 3)$
- How to calculate E[V]?

TWO-DIMENSIONAL CHANGE-OF-VARIABLE FORMULA. Let X and Y be random variables, and let $g: \mathbb{R}^2 \to \mathbb{R}$ be a function. If X and Y are *discrete* random variables with values a_1, a_2, \ldots and b_1, b_2, \ldots , respectively, then

$$\mathbf{E}[g(X,Y)] = \sum_{i} \sum_{j} g(a_i, b_j) \mathbf{P}(X = a_i, Y = b_j) \,.$$

If X and Y are continuous random variables with joint probability density function f, then

$$\mathrm{E}\left[g(X,Y)\right] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) f(x,y) \,\mathrm{d}x \,\mathrm{d}y.$$

If $H \perp \!\!\!\perp R$:

$$E[V] = E[\pi HR^2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \pi hr^2 f_H(h) f_R(r) dh dr$$

Linearity of expectations

Theorem. For X and Y random variables, and $s, t \in \mathbb{R}$:

$$E[rX + sY + t] = rE[X] + sE[Y] + t$$

Proof. (discrete case)

$$E[rX + Ys + t] = \sum_{a} \sum_{b} (ra + sb + t)P(X = a, Y = b)$$

= $\left(r\sum_{a} \sum_{b} aP(X = a, Y = b)\right) + \left(s\sum_{a} \sum_{b} bP(X = a, Y = b)\right) + \left(t\sum_{a} \sum_{b} P(X = a, Y = b)\right)$
= $\left(r\sum_{a} aP(X = a)\right) + \left(s\sum_{b} bP(Y = b)\right) + t = rE[X] + sE[Y] + t$

Corollary. $E[a_0 + \sum_{i=1}^n a_i X_i] = a_o + \sum_{i=1}^n a_i E[X_i]$

Corollary. $X \le Y$ implies $E[X] \le E[Y]$ **Proof.** $Z = Y - X \ge 0$ implies $E[Z] = E[Y] - E[X] \ge 0$, i.e., $E[Y] \ge E[X]$.

Applications

- Expectation of some discrete distributions
 - $X \sim Ber(p)$ E[X] = p
 - $X \sim Bin(n, p)$ $E[X] = n \cdot p$
 - $\square \text{ Because } X = \sum_{i=1}^{n} X_i \text{ for } X_1, \dots, X_n \sim Ber(p)$
 - $X \sim Geo(p)$ $E[X] = \frac{1}{p}$

►
$$X \sim NBin(n, p)$$
 $E[X] = \frac{n \cdot (1-p)}{p}$
□ Because $X = \sum_{i=1}^{n} X_i - n$ for $X_1, \dots, X_n \sim Geo(p)$

• Expectation of some continuous distributions

$$\begin{array}{l} \blacktriangleright X \sim Exp(\lambda) & E[X] = \frac{1}{\lambda} \\ \blacktriangleright X \sim Erl(n,\lambda) & E[X] = \frac{n}{\lambda} \\ \Box \text{ Because } X = \sum_{i=1}^{n} X_i \text{ for } X_1, \dots, X_n \sim Exp(\lambda) \end{array}$$

Expectation of product and quotients

Theorem. For $X \perp Y$, we have: E[XY] = E[X]E[Y]

PROPAGATION OF INDEPENDENCE. Let X_1, X_2, \ldots, X_n be independent random variables. For each i, let $h_i : \mathbb{R} \to \mathbb{R}$ be a function and define the random variable

 $Y_i = h_i(X_i).$

Then Y_1, Y_2, \ldots, Y_n are also independent.

Corollary. For $X \perp Y$ and Y > 0, we have: $E[X/Y] \ge E[X]/E[Y]$ *Proof.* $X \perp Y$ implies $X \perp 1/Y$. By theorem above:

 $E[X/Y] = E[X \cdot 1/Y] = E[X]E[1/Y] \ge E[X]/E[Y]$

because by Jensen's inequality $E[1/r] \ge 1/E[Y]$ since 1/y is convex for y = 0. **Exercise at home.** Show that E[X/Y] = E[X]/E[Y] is a false claim.

Prove it!

Law of iterated/total expectation

Conditional expectation

$$E[X|Y = b] = \sum_{i} a_{i}p(a_{i}|b) \qquad E[X|Y = y] = \int_{-\infty}^{\infty} xf(x|y)dx$$

Theorem. (Law of iterated/total expectation)

$$E_Y[E[X|Y]] = E[X]$$

Proof. (for X, Y discrete random variables)

$$E_{Y}[E[X|Y]] = \sum_{j} \sum_{i} a_{i} p_{X|Y}(a_{i}|b_{j}) p_{Y}(b_{j}) = \sum_{j} \sum_{i} a_{i} p_{XY}(a_{i}, b_{j}) = \sum_{i} a_{i} p_{X}(a_{i}) = E[X]$$

Example (cfr the example from Lesson 1 on the Law of total probability)

- Factory 1's light bulbs working hours $\sim \textit{Exp}(1/1000)$
- Factory 2's light bulbs working hours $\sim Exp(1/2000)$
- Factory 1 supplies 60% of the total bulbs on the market and Factory 2 supplies 40% of it.
- What is the average work hour of a light bulb on the market?

Variance of the sum and covariance

 $Var(X + Y) = E[(X + Y - E[X + Y])^{2}] = E[((X - E[X]) + (Y - E[Y]))^{2}]$

$$= E[(X - E[X])^{2}] + E[(Y - E[Y])^{2}] + 2E[(X - E[X])(Y - E[Y])]$$

$$=$$
 Var(X) + Var(Y) + 2Cov(X, Y)

Covariance

The covariance Cov(X, Y) of two random variables X and Y is the number:

Cov(X, Y) = E[(X - E[X])(Y - E[Y])]



Covariance

Theorem. Cov(X, Y) = E[XY] - E[X]E[Y]

• If X and Y are independent $(X \perp \!\!\!\perp Y)$:

$$Cov(X, Y) = 0$$
 $Var(X + Y) = Var(X) + Var(Y)$

- But there are X and Y uncorrelated (ie., Cov(X, Y) = 0) that are dependent!
- Variances of some discrete distributions
 - $X \sim Ber(p)$ Var(X) = p(1-p)
 - ► $X \sim Bin(n, p)$ Var(X) = np(1-p)□ Because $X = \sum_{i=1}^{n} X_i$ for $X_1, \dots, X_n \sim Ber(p)$ and independent
 - $X \sim Geo(p)$ $Var(X) = \frac{1-p}{p^2}$
 - $X \sim NBin(n, p)$ $Var(X) = n \frac{1-p}{p^2}$

 \square Because $X = \sum_{i=1}^{n} X_i - n$ for $X_1, \ldots, X_n \sim Geo(p)$ and independent

- Variances of some continuous distributions
 - $X \sim Exp(\lambda)$ $Var(X) = \frac{1}{\lambda^2}$ • $X \sim Frl(n, \lambda)$ $Var(X) = \frac{n}{2}$

$$\square \text{ Because } X = \sum_{i=1}^{n} X_i \text{ for } X_1, \dots, X_n \sim Exp(\lambda) \text{ and independent}$$

COVARIANCE UNDER CHANGE OF UNITS. Let X and Y be two random variables. Then

```
\operatorname{Cov}(rX + s, tY + u) = rt\operatorname{Cov}(X, Y)
```

for all numbers r, s, t, and u.

- Hence, $Var(rX + sY + t) = r^2 Var(X) + s^2 Var(Y) + 2rsCov(X, Y)$
- **Bivariate** Normal/Gaussian distribution:

$$(X, Y) \sim \mathcal{N}((\mu_X, \mu_X), \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{pmatrix})$$

- where marginals are $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$, $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$, and $Cov(X, Y) = \sigma_{XY}$
- Covariance matrix $\Sigma_{ij} = Cov(X_i, X_j)$ for a vector $\mathbf{X} = (X_1, \dots, X_n)$ of r.v.'s
- Covariance depends on the unit of measure!

DEFINITION. Let X and Y be two random variables. The correlation coefficient $\rho(X, Y)$ is defined to be 0 if $\operatorname{Var}(X) = 0$ or $\operatorname{Var}(Y) = 0$, and otherwise $\rho(X, Y) = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}}.$

- Correlation coefficient is *dimensionless* (not affected by change of units)
 - E.g., if X and Y are in Km, then Cov(X, Y), Var(X) and Var(Y) are in Km²
- Moreover: $-1 \leq
 ho(X,Y) \leq 1$
 - ► The bounds are derived from the Cauchy-Schwarz's inequality:

$$E[|XY|] \le \sqrt{E[X^2]} \sqrt{E[Y^2]}$$

Proof. For any $u, w \in \mathbb{R}$, we have $2|uw| \le u^2 + w^2$. Therefore, $2|UW| \le U^2 + W^2$ for r.v.'s U and V. By defining $U = \frac{x}{\sqrt{E[x^2]}}$ and $W = \frac{Y}{\sqrt{E[Y^2]}}$ (*), we have $2 \cdot \frac{|XY|}{\sqrt{E[X^2]}\sqrt{E[Y^2]}} \le \frac{x^2}{E[X^2]} + \frac{Y^2}{E[Y^2]}$. Taking the expectations, we conclude: $2 \cdot \frac{E[|XY|]}{\sqrt{E[X^2]}\sqrt{E[Y^2]}} \le 2$. (*) The case $E[X^2] = 0$ or $E[Y^2] = 0$ is left as an exercise.

Bivariate Normal/Gaussian distribution

$$(X, Y) \sim \mathcal{N}((\mu_X, \mu_Y), \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{pmatrix})$$

where marginals are $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$, $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$, and $Cov(X, Y) = \sigma_{XY}$

• Since
$$\sigma_{XY} = \rho(X, Y) \cdot \sigma_X \cdot \sigma_Y$$
:
 $(X, Y) \sim \mathcal{N}((\mu_X, \mu_Y), \begin{pmatrix} \sigma_X^2 & \rho(X, Y) \cdot \sigma_X \cdot \sigma_Y \\ \rho(X, Y) \cdot \sigma_X \cdot \sigma_Y & \sigma_Y^2 \end{pmatrix})$

• Density of $\mathcal{N}((0, 0), (1, \sigma_{XY}, \sigma_{XY}, 1))$:

$$f(x,y) = \frac{1}{2\pi\sqrt{1-\sigma_{XY}^2}}e^{-\frac{1}{2(1-\sigma_{XY}^2)}(x^2+y^2-2xy\sigma_{XY})}$$

- Useful facts for (X, Y) bivariate Normal:
 - ▶ for (X, Y) bivariate Normal: $\rho(X, Y) = 0$ iff $X \perp Y$, i.e., uncorrelation equals independence
 - ▶ (X, Y) bivariate Normal iff aX + bY is Normal for any $a, b \in \mathbb{R}$

Sum of independent Normal random variables

• See Lesson 04 and Lesson 08 for convolution formulas

ADDING TWO INDEPENDENT CONTINUOUS RANDOM VARIABLES. Let X and Y be two independent continuous random variables, with probability density functions f_X and f_Y . Then the probability density function f_Z of Z = X + Y is given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(z-y) f_Y(y) \,\mathrm{d}y$$

for
$$-\infty < z < \infty$$
.

Theorem. If $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ and $X \perp Y$, then: $Z = X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$

Proof. See [T, Sect. 11.2]

• In general: $Z = rX + sY + t \sim \mathcal{N}(r\mu_X + s\mu_Y + t, r^2\sigma_X^2 + s^2\sigma_Y^2)$

The converse of the theorem also holds:

[Lévy-Cramér theorem]

• If $X \perp Y$ and Z = X + Y is normally distributed, then X and Y follow a normal distribution.

Extremes of independent random variables

THE DISTRIBUTION OF THE MAXIMUM. Let X_1, X_2, \ldots, X_n be n independent random variables with the same distribution function F, and let $Z = \max\{X_1, X_2, \ldots, X_n\}$. Then

 $F_Z(a) = (F(a))^n.$

•
$$P(Z \le a) = P(X_1 \le a, ..., X_n \le a) = \prod_{i=1}^n P(X_i \le a) = ((F(a))^n)$$

- Example: maximum water level over 365 days assuming water level on a day is U(0,1)
- Example: maximum of two rolls of a die with 4 sides

THE DISTRIBUTION OF THE MINIMUM. Let X_1, X_2, \ldots, X_n be n independent random variables with the same distribution function F, and let $V = \min\{X_1, X_2, \ldots, X_n\}$. Then

$$F_V(a) = 1 - (1 - F(a))^n.$$

•
$$P(V \le a) = 1 - P(X_1 > a, ..., X_n > a) = 1 - \prod_{i=1}^n (1 - P(X_i \le a)) = 1 - ((1 - F(a))^n)$$

Product and quotient of independent random variables

PRODUCT OF INDEPENDENT CONTINUOUS RANDOM VARIABLES. Let X and Y be two independent continuous random variables with probability densities f_X and f_Y . Then the probability density function f_Z of Z = XY is given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_Y\left(\frac{z}{x}\right) f_X(x) \frac{1}{|x|} \, \mathrm{d}x$$

for $-\infty < z < \infty$.

QUOTIENT OF INDEPENDENT CONTINUOUS RANDOM VARIABLES. Let X and Y be two independent continuous random variables with probability densities f_X and f_Y . Then the probability density function f_Z of Z = X/Y is given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(zx) f_Y(x) |x| \, \mathrm{d}x$$

for $-\infty < z < \infty$.

• $X, Y \sim \mathcal{N}(0,1)$ independent, $Z = X/Y \sim Cau(0,1)$ where:

$$f_Z(x) = \frac{1}{\pi(1+x^2)}$$

Distances and Metrics

A numerical measurement of how far apart two objects are.

Distances and Metrics

A distance over a set \mathcal{A} is a function $d : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$ such that:

- $d(x,y) \ge 0$ non-negativity
- d(x, y) = 0 iff x = y
- d(x,y) = d(y,x)

Moreover, d is called a metric if in addition:

• $d(x,z) \leq d(x,y) + d(y,z)$

identity of indiscernibles

symmetry

triangle inequality

Examples over $\mathcal{A} = \mathbb{R}^n$:

- Manhattan or L_1 distance $d_1(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} \mathbf{y}\|_1 = \sum_{i=1}^n |\mathbf{x}_i \mathbf{y}_i|$
- Euclidian or L_2 distance $d_2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} \mathbf{y}\|_2 = \sqrt{\sum_{i=1}^n (\mathbf{x}_i \mathbf{y}_i)^2}$

• Chebyshev or L_∞ distance $d_\infty(\mathbf{x},\mathbf{y}) = \|\mathbf{x}-\mathbf{y}\|_\infty = \max_{i=1}^n |\mathbf{x}_i-\mathbf{y}_i|$

We aim at defining distances and metrics over probability distributions, i.e., when $\mathcal{A} = \{F \mid F : \mathbb{R} \to [0, 1] \text{ is a CDF} \}$

Distances over probability distributions

A numerical measurement of how far apart two probability distributions are.

- ML/DM models are supposed to be applied on the same distribution as the training set:
 - How far is the test data distribution from the one of the training data?
 - Is the data changing over time, thus my model is inadequate?
- ML/DM algorithms are supposed to choose the best hypothesis:
 - What is the split in a DT which best distinguish the distribution of classes?
 - Is my model separating positive and negatives as much as possible?
 - Is my clustering separating groups with different distributions?
- Data preprocessing looks at feature distribution:
 - Are these two features conveying the same information?
 - Can this feature be predictive to the class feature?
- ... and many other applications in Data Science

[Transfer learning [Dataset shift]

Total variation distance and KS distance

- Let X, Y be random variables:
 - Total Variation (TV) distance (discrete and continuous case):

$$d_{TV}(X,Y) = \frac{1}{2} \sum_{i} |p_X(a_i) - p_Y(a_i)| \qquad d_{TV}(X,Y) = \frac{1}{2} \int |f_X(x) - f_Y(x)| dx$$

- d_{TV} is a metric with $d_{TV}(X, Y) \in [0, 1]$
- Kolmogorov-Smirnov (KS) distance:

$$d_{\mathcal{KS}}(X,Y) = \sup_{x} |F_X(x) - F_Y(x)|$$

- d_{KS} is a metric with $d_{KS}(X, Y) \in [0, 1]$
- d_{TV} and d_{KS} have no closed forms in general
- d_{KS} can be estimated from samples of the distributions



Entropy H(X) of a random variable X

- The Shannon's information entropy is the average level of "information" (or "surprise", "uncertainty", "unpredictability") inherent to the variable's possible outcomes
 - Information is inversely proportional to probability
 - $\hfill\square$ Highly likely/unlikely events carry less/more new information
 - Information content ic() of two independent events should sum up
 - $\begin{array}{l} \square \quad ic(p(A \cap B)) = ic(p(A)) + ic(p(B)) = ic(p(A)p(B)) \\ \square \quad ic(p(\Omega)) = ic(1) = 0 \\ \square \quad ic(p(A)) \ge 0 \end{array}$

•
$$H(X) = E[-\log p(X)]$$
 (discrete)
 $H(X) = E[-\log f(X)]$ (continuous)
 $H(X) = -\sum_{i} p(a_i) \log p(a_i)$
 $H(X) = -\int_{-\infty}^{\infty} f(x) \log f(x) dx$

▶ For X discrete, $H(X) \ge 0$ since $-\log p(X) = \log \frac{1}{p(X)} \ge 0$ □ zero reached when $p(a_1) = 1$ and $p(a_i) = 0$ for $i \ne 1$

► For
$$X \sim Ber(p)$$
, $H(X) = -p \log p - (1-p) \log (1-p)$
□ for $X \sim Ber(0.5)$: $H(X) = -2 \cdot \frac{1}{2} \log \frac{1}{2} = 1$

[binary entropy function] [unit of entropy is called a bit]

 $\frac{1}{p(a)}$

 $\log \frac{1}{p(a_i)}$

Cross entropy

- X, Y discrete random variables with p.m.f. p_X and p_Y :
- Cross entropy of X w.r.t. Y: $H(X; Y) = E_X[-\log p(Y)]$

$$H(X; Y) = -\sum_{i} p_X(a_i) \log p_Y(a_i)$$

with $p_X(a_i) \log p_Y(a_i) = \begin{cases} 0 & \text{if } p_X(a_i) = 0 \\ -\infty & \text{if } p_X(a_i) > 0 \land p_Y(a_i) = 0 \end{cases}$

- H(X; Y) is the "information" or "uncertainty" or "loss" when using Y to encode X
- The closer p_X and p_Y , the lower is H(X; Y)
- The lower bound is for Y = X, for which H(X; Y) = H(X)

Kullback-Leibler divergence

KL divergence

For X, Y discrete random variables with p.m.f. p_X and p_Y :

$$D_{\mathcal{KL}}(X \parallel Y) = \sum_i p_X(a_i) \log \frac{p_X(a_i)}{p_Y(a_i)} = H(X;Y) - H(X)$$

- Measure how distribution of Y (model) can reconstruct the distribution of X (data)
 - ► Also called: relative entropy or information gain of X w.r.t. Y
- Properties
 - $D_{KL}(X \parallel Y) \geq 0$
 - $D_{KL}(X \parallel Y) = 0$ iff $F_X = F_Y$
 - $D_{KL}(X \parallel Y) \neq D_{KL}(Y \parallel X)$
- For X, Y continuous: $D_{KL}(X \parallel Y) = \int_{-\infty}^{\infty} f_X(x) \log \frac{f_X(x)}{f_Y(x)} dx$

[Gibbs' inequality]

[not a distance!]

Joint entropy

- X, Y discrete random variables with p.m.f. p_X and p_Y :
- Joint p.m.f. p_{XY} . Joint entropy of (X, Y):

$$H((X,Y)) = -\sum_{i,j} p_{XY}(a_i,a_j) \log p_{XY}(a_i,a_j)$$

• If $X \perp Y$, then:

$$H((X, Y)) = -\sum_{i,j} p_X(a_i) p_Y(a_j) (\log p_X(a_i) + \log p_Y(a_j)) =$$

= -(\sum_i p_X(a_i)) (\sum_j p_Y(a_j) \log p_Y(a_j)) - (\sum_j p_Y(a_j)) (\sum_i p_X(a_i) \log p_X(a_i)) = H(X) + H(Y)

Mutual information

Mutual information

For X, Y discrete random variables with p.m.f. p_X and p_Y and joint p.m.f. p_{XY} :

$$I(X, Y) = D_{KL}(p_{XY} \parallel p_X p_Y) = \sum_{i,j} p_{XY}(a_i, a_j) \log \frac{p_{XY}(a_i, a_j)}{p_X(a_i)p_Y(a_j)} = H(X) + H(Y) - H((X, Y))$$

- MI measures how dependent two distributions are
 - Measure how product of marginals can reconstruct the joint distribution
- Properties

•
$$I(X,Y) = I(Y,X)$$
, and $I(X,Y) \ge 0$

•
$$I(X, Y) = 0$$
 iff $X \perp Y$

•
$$NMI = \frac{I(X,Y)}{\min \{H(X),H(Y)\}} \in [0,1]$$

[Normalized mutual information]

• For X, Y continuous: $I(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) \log \frac{f_{XY}(x, y)}{f_X(x)f_Y(y)} dxdy$

The data processing inequality

- Let X be unknown, and assume to observe a noisy version Y of it
- Let Z = f(Y) be a data processing to improve the "quality" of Y
- Z does not increase the information about X, i.e.:

 $I(X, Y) \geq I(X, Z)$

- If I(X, Y) = I(X, Z) and Z is a summary of Y, we call it a sufficient statistics
 - Let $X \sim Ber(\theta)$ and $Y = (Y_1, \ldots, Y_n) \sim Ber(\theta)^n$ modelling i.i.d. observations
 - $Z = \sum_{i=1}^{n} Y_i \sim Binom(n, \theta)$ is a sufficient statistics
 - **Proof (sketch):** use $D_{KL}(p_{XY} \parallel p_X p_Y)$ and:

$$p(Y_1 = y_1, ..., Y_n = y_n) = \prod_i \theta^{y_i} (1 - \theta)^{(1 - y_i)} = \theta^{\sum_i y_i} (1 - \theta)^{n - \sum_i y_i} = p(Z = \sum_i y_i)$$

[Data processing inequality]

Earth mover's distance / Wasserstein metric

- The minimum cost to transform one distribution to another
- Cost = amount of mass to move \times distance to move it
- X, Y discrete random variables:

$$\mathsf{EMD}(X,Y) = rac{\sum_{i,j} \mathsf{F}_{i,j} \cdot |\mathsf{a}_i - \mathsf{a}_j|}{\sum_{i,j} \mathsf{F}_{i,j}}$$

where F is the flow which minimizes the numerator (total cost) subject to some constraints.



Earth mover's distance / Wasserstein metric

- The minimum cost to transform one distribution to another
- Solution of the transportation problem for X, Y multivariate (version from Ramdas et al. 2015):

$$EMD(X, Y) = \int_0^1 \|F_X^{-1}(p) - F_Y^{-1}(p)\| dp$$

For X, Y univariate, this simplifies to:

$$EMD(X,Y) = \sum_i |F_X(a_i) - F_Y(a_i)|$$
 $EMD(X,Y) = \int_{-\infty}^{\infty} |F_X(x) - F_Y(x)| dx$

• For empirical distributions from **ordered** samples x_1, \ldots, x_n and y_1, \ldots, y_n :

$$EMD(X, Y) = \frac{1}{n}\sum_{i}|x_i - y_i|$$

Simulation

- Not all problems can be solved with calculus!
- Complex interactions among random variables can be simulated
- Generated random values are called *realizations*
- Basic issue: how to generate realizations?
 - ► The Galton Board



Simulation

- Not all problems can be solved with calculus!
- Complex interactions among random variables can be simulated
- Generated random values are called *realizations*
- Basic issue: how to generate realizations?
- Assumption: we are only given U(0,1)
- Problem: derive all the other random generators

Simulation: discrete distributions

Bernoulli random variables

Suppose U has a U(0, 1) distribution. To construct a Ber(p) random variable for some 0 , we define

$$X = \begin{cases} 1 & \text{if } U < p, \\ 0 & \text{if } U \ge p \end{cases}$$

so that

$$P(X = 1) = P(U < p) = p,$$

 $P(X = 0) = P(U \ge p) = 1 - p.$

This random variable X has a Bernoulli distribution with parameter p.

• For $X_1, \ldots, X_n \sim Ber(p)$ i.i.d., we have: $\sum_{i=1}^n X_i \sim Binom(n, p)$

$X \sim Cat(\mathbf{p})$

DEFINITION. A discrete random variable X has a *Bernoulli distribution* with parameter p, where $0 \le p \le 1$, if its probability mass function is given by

 $p_X(1) = P(X = 1) = p$ and $p_X(0) = P(X = 0) = 1 - p$.

We denote this distribution by Ber(p).

- Alternative definition: $p_X(a) = p^a \cdot (1-p)^{1-a}$ for $a \in \{0,1\}$
- Categorical distribution generalizes to $n_C \ge 2$ possible values

$$X \sim Cat(\mathbf{p})$$

Categorical distribution

A discrete random variable X has a Categorical distribution with parameters p_0, \ldots, p_{n_c-1} where $\sum_i p_i = 1$ and $p_i \in [0, 1]$ if its p.m.f. is given by:

 $p_X(i) = P(X = i) = p_i$ for $i = 0, ..., n_C - 1$

• Alternative definition: $p_X(a) = \prod_i p_i^{\mathbb{1}_{a=i}}$ for $a \in \{0, \dots, n_C - 1\}$

Notation. Indicator function: $\mathbb{1}_{\varphi}(x) = \begin{cases} 1 & \text{if } \varphi(x) \\ 0 & \text{otherwise} \end{cases}$

$X \sim Mult(n, \mathbf{p})$

- $X \sim Bin(n, p)$ models the number of successes in *n* Bernoulli trials
- Intuition: for X_1, X_2, \ldots, X_n i.i.d. $X_i \sim Ber(p)$: $X = \sum_{i=1}^n X_i \sim Bin(n, p)$
- X ~ Mult(n, **p**) models the number of categories in n Categorical trials
- Intuition: for X_1, X_2, \ldots, X_n such that $X_i \sim Cat(\mathbf{p})$ and independent (i.i.d.), define:

$$Y_{1} = \sum_{i=1}^{n} \mathbb{1}_{X_{i}=0} \sim Bin(n, p_{0}) \quad \dots \quad Y_{n_{c}} = \sum_{i=1}^{n} \mathbb{1}_{X_{i}=n_{c}-1} \sim Bin(n, p_{n_{c}-1})$$
$$X = (Y_{1}, \dots, Y_{n_{c}}) \sim Mult(n, \mathbf{p})$$

Multinomial distribution

A discrete random variable $X = (Y_1, \ldots, Y_{n_c})$ has a Multinomial distribution with parameters p_0, \ldots, p_{n_c-1} where $\sum_i p_i = 1$ and $p_i \in [0, 1]$ if its p.m.f. is given by:

$$p_X(i_0,\ldots,i_{n_c-1}) = P(X = (i_0,\ldots,i_{n_c-1})) = \frac{n!}{i_0!i_1!\ldots i_{n_c-1}!} p_0^{i_0} p_1^{i_1} \ldots p_{(n_c-1)}^{i_{(n_c-1)}}$$

$X \sim Mult(n, \mathbf{p})$

- Example: student selection from a population with $n_C = 3$:
 - $p_0 = 60\%$ undergraduates
 - $p_1 = 30\%$ graduate
 - $p_2 = 10\%$ PhD students
- Assume n = 20 students are randomly selected
- $X \sim (Y_1, Y_2, Y_3)$ where:
 - ► Y₁ number of undergraduate students selected
 - ► Y₂ number of graduate students selected
 - ► Y₃ number of PhD students selected
- $P(X = (10, 6, 4)) = \frac{20!}{10!6!4!} (0.6)^{10} (0.3)^6 (0.1)^4 = 9.6\%$

Simulation: continuous distributions

- $F(x) = P_X(X \leq x)$
- $F:\mathbb{R} \to [0,1]$ invertible as $F^{-1}:[0,1] \to \mathbb{R}$
 - ► E.g., *F* strictly increasing
 - N.B., the textbook notation for F^{-1} is F^{inv}
- For $Y \sim U(0,1)$ and $0 \le b \le 1$ $P_Y(Y \le b) = b$

then, for b = F(x) $P_Y(Y \le F(x)) = F(x)$

and then by inverting $X = F^{-1}(Y)$ $P_X(X \le x) = P_Y(F^{-1}(Y) \le x) = F(x)$

• In summary:

$$X=F^{-1}(Y)\sim F$$
 for $Y\sim U(0,1)$



Kevin P. Murphy (2022)

Probabilistic Machine Learning: An Introduction Chapter 6: Information Theory online book

William H. Press, Saul A. Teukolsky, William T. Vetterling, Brian P. Flannery (2007) Numerical Recipes - The Art of Scientific Computing Chapter 7: Random Numbers online book