# Probabilistic Generative models
## Week 1 — Foundations

A PGM, or simply a GM, is a methodology for generating data. In general, the GM is constructed and adjusted using observations with the aim to generate data with the same statistical properties of the observations. The 'probabilistic' nature of the PGMs considered follow from the fact that the data generated will be considered to be realisations of a Random variable; e.g. $X$.

In this sense, the probability distribution of $X$, denoted $P_X(x)$, as well as its density $p_X(x)$ will be central to the study of PGM. In particular, targeting the pdf is one way of constructing PGMs, in which case the whole PGM paradigm becomes equivalent to the classical statistical model approach. However, as we will see in the course, enforcing the sought-after PGM to have an explicit parametric pdf can be rather restrictive.

Throughout the course, we will consider a probability space $(\Omega, \mathcal{F}, P)$ with 3 RVs given by the measurable maps

$$X: \Omega \to \mathcal{X} \quad ; \quad Y: \Omega \to \mathcal{Y} \quad ; \quad Z: \Omega \to \mathcal{Z}$$

(observed input)    (observed output)    (latent variable)

Remark: Not all three variables will be present in all our models. For instance, in classification there's no justification for $Z$, while in clustering there's no need for $Y$. However, we build the general set up here for formality.

We will also consider the $\sigma$-algebra $\mathcal{F}$ to be the product Borel $\sigma$-algebra. $\mathcal{F} = \mathcal{B}(X) \otimes \mathcal{B}(Y) \otimes \mathcal{B}(Z)$.

Furthermore, we will assume the joint probability of $X$, $Y$ and $Z$ has a density, that is, $\forall A \in \mathcal{B}(X)$, $B \in \mathcal{B}(Y)$, $C \in \mathcal{B}(Z)$, we have

$$P(X \in A, Y \in B, Z \in C) = \int_{A \times B \times C} p(x,y,z) \, dx \, dy \, dz$$

We will also assume that all marginals and conditionals, have a density. This includes $p(x,y)$, $p(y|x)$, $p(z|x,y)$ etc.

Discriminative $V/s$ Generative. The generative approach aims to characterise the complete generative distribution $p(x,y,z)$, whereas in some application specific cases only the discriminative model $p(y|x)$ is needed. Let's a-the following example.

Example: generative $V/s$ discriminative $\overset{\text{binary}}{\text{classification}}$.

Consider the $\overset{\text{bi}}{\text{classification}}$ problem, where given an observation $X = x$, one needs to estimate its label $Y$. A discriminative model would directly construct $p(y|X=x)$. Since this is a binary case, without loss of generality we can assume $Y \in \{0,1\}$, and model $P(Y=1|X=x)$, since $P(Y=0|X=x) = 1 - \%$ A model for this will only need to map $x \in \mathbb{R}^u \rightarrow P(Y=1|X=x) \in [0,1]$. For instance

$$P(Y=1|X=x) = \frac{1}{1 + e^{-\sigma^T x}}$$

which is known as the logistic regression.

Conversely, in a generative approach, we aim to model the joint probability $p(Y=1, X=x)$.

Modelling this distribution directly is not easy, however, we can observe that it can be factorised as

$$p(Y=1, X=x) = p(X=x|Y=1)\, p(Y=1)$$

which yields two much more intuitive distributions:

- the class prob $p(Y=1) \to \pi,\ (1-\pi)$
- the class-conditional prob $p(X=x|Y=1) \to$ choose...
  $$f_{\theta_1}, f_{\theta_2}$$

Therefore, the classifier is

$$p(Y=1|X=x) = \frac{p(X=x|Y=1)\, p(Y=1)}{p(X=x)}$$

$$= \frac{1}{1 + \frac{p(X=x|Y=0)p(Y=0)}{p(X=x|Y=1)p(Y=1)}}$$

$$= \frac{1}{1 + e^{\log\left(\frac{1-\pi}{\pi} \cdot \frac{f_{\theta_0}(X)}{f_{\theta_1}(X)}\right)}} \qquad \circledast$$

Exercise: evaluate $\circledast$ for $f_{\theta_1} = N(\mu_1, \Sigma),\ f_{\theta_0} = N(\mu_2, \Sigma)$

- MLE example
- push forwards
- iid theo
- pdf not always possible