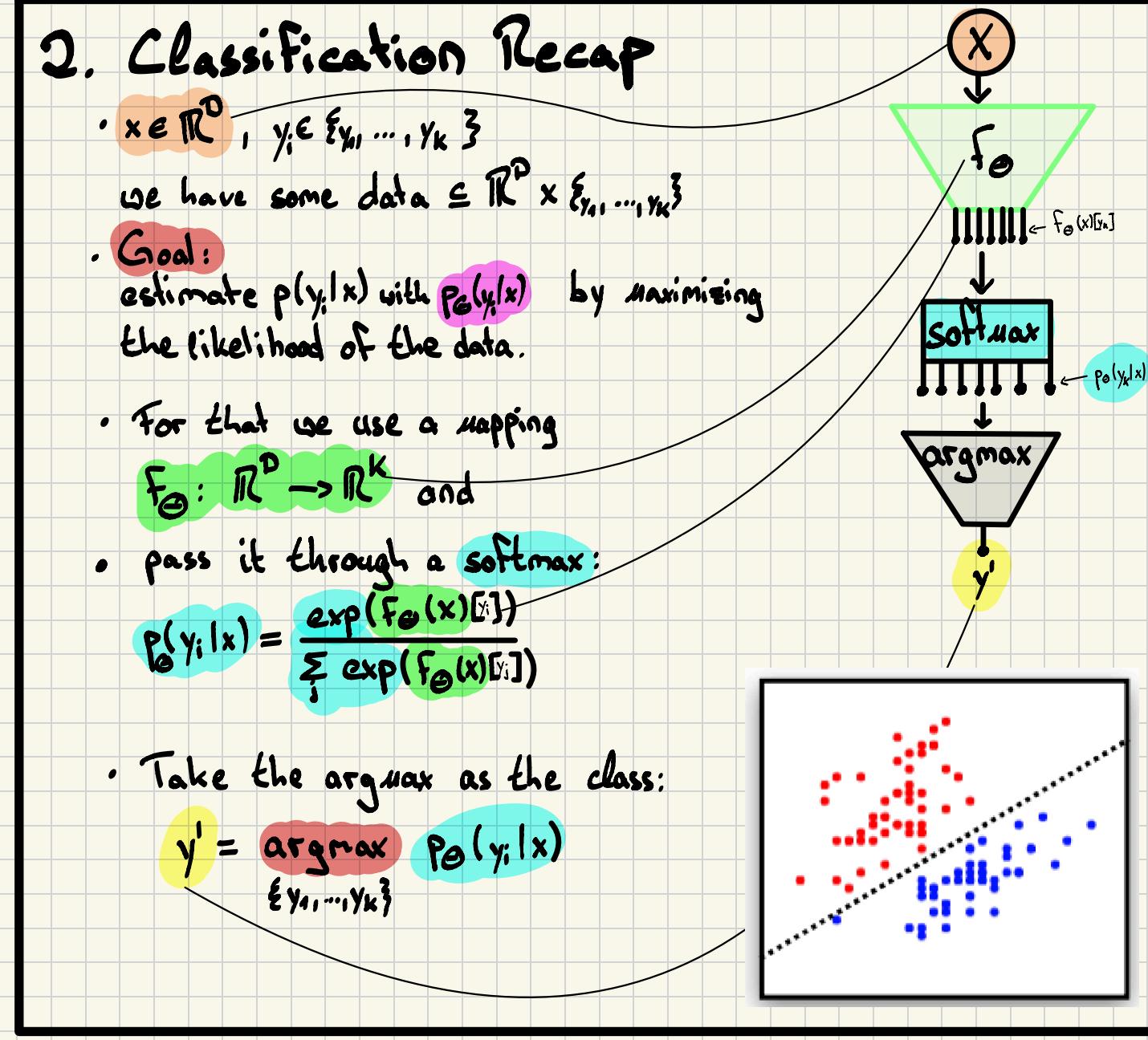
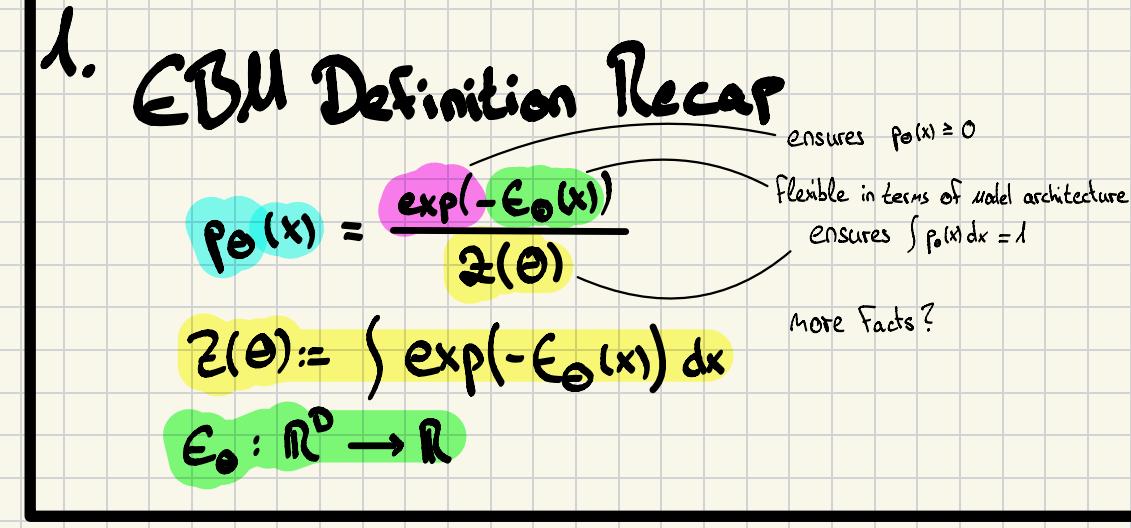


Your classifier is secretly an EBM and you should treat it like one



- 3. A different perspective:**
 - Redefine the meaning of embeddings

• EBM: $p_\theta(x, y_i) = \frac{\exp(F_\theta(x)[y_i])}{Z(\theta)}$

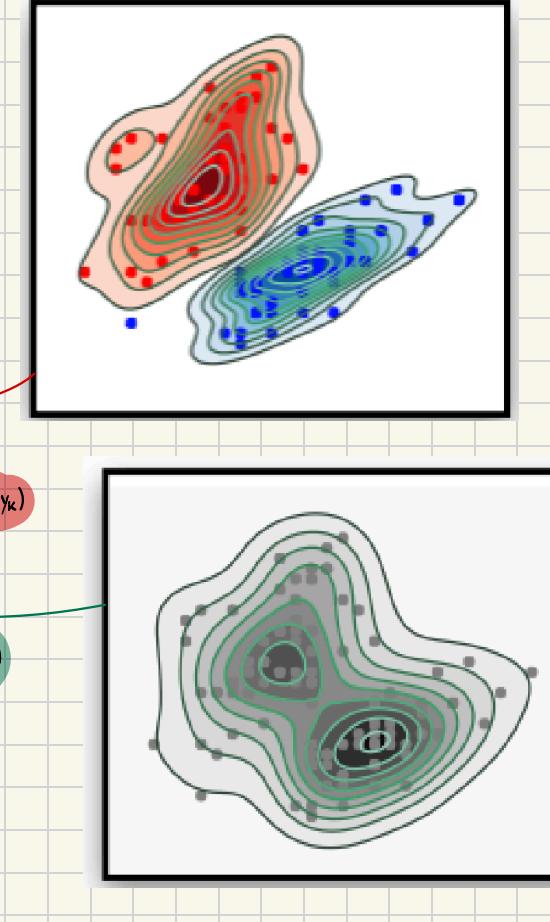
• with Energy $E_\theta(x, y_i) = -F_\theta(x)[y_i]$

• Aggregation over classes gives us the marginal: an EBM for $p(x)$:

$$p_\theta(x) = \sum_{y \in Y} p_\theta(x, y) = \frac{1}{Z(\theta)} \sum_{y \in Y} \exp(F_\theta(x)[y])$$
 with energy $E_\theta(x) = -\log \sum_{y \in Y} \exp(F_\theta(x)[y])$

• Conditioning on x recovers the softmax!

$$p_\theta(y_i|x) = \frac{p_\theta(x, y_i)}{p_\theta(x)} = \frac{\frac{\exp(F_\theta(x)[y_i])}{Z(\theta)}}{\frac{1}{Z(\theta)} \sum_{y \in Y} \exp(F_\theta(x)[y])} = \frac{\exp(F_\theta(x)[y_i])}{\sum_{y \in Y} \exp(F_\theta(x)[y])}$$



- 4. Insight: We have 3 EBMs hidden in the classifier.**

1. The obvious conditional $p_\theta(y_i|x)$ (actually K conditionals and joint models)
2. The joint $p_\theta(x, y_i)$
3. The marginal $p_\theta(x)$

5. How 2 train?

$$\frac{p_\theta(x, y_i)}{p_\theta(x)} = p_\theta(y_i|x)$$

$$\log p_\theta(x, y_i) - \log p_\theta(x) = \log p_\theta(y_i|x)$$

$$\text{Loss}(x, y_i, \theta) = -\log p_\theta(x, y_i) = -(\log p_\theta(y_i|x) + \log p_\theta(x))$$

Vanilla Cross-Entropy loss
Estimate using Stochastic Gradient Langevin Dynamics

6. Stochastic Gradient Langevin Dynamics (SGLD)

Recap:

$$p_\theta(x) = \sum_{y \in Y} p_\theta(x, y) = \frac{1}{Z(\theta)} \sum_{y \in Y} \exp(F_\theta(x)[y])$$

$$E_\theta(x) = -\log \sum_{y \in Y} \exp(F_\theta(x)[y])$$

Derivative of the log-likelihood for a single example x with resp. to the parameters θ can be expressed as:

$$\frac{\partial \log p_\theta(x)}{\partial \theta} = \mathbb{E}_{p_\theta(x)} \left[\frac{\partial \exp(F_\theta(x))}{\partial \theta} \right] - \frac{\partial \exp(F_\theta(x))}{\partial \theta}$$

Intractable
 → MCMC-Sampling with SGLD + Persistent Contrastive (PCD)
 Divergence for $p_\theta(x)$ for faster sampling convergence

- SGLD**
- (1) $x \sim p_\theta(x)$ (which is typically a uniform dist over the input domain)
 - (2) $\alpha_t \in \mathbb{R}^+$ (which should follow a polynomial decay schedule)
 - (3) for some iterations do:
 - (3.1) $\epsilon_t \sim N(0, \alpha_t I)$
 - (3.2) $x_{t+1} = x_t - \frac{\alpha_t}{2} \frac{\partial \exp(F_\theta(x_t))}{\partial x_t} + \epsilon_t$
 - (3.3) $\alpha_{t+1} = \text{PolynomialDecaySchedule}(\alpha_t)$

PCD
 but with PCD consists of a fraction $B \in (0, 1)$ of the samples from the last step and $(1-B)$ fraction of uniform noise.

Implementation Demo

EBM Definition Recap

$$p_\theta(x) = \frac{\exp(-E_\theta(x))}{Z(\theta)}$$

$$Z(\theta) = \int \exp(-E_\theta(x)) dx$$

$$E_\theta: \mathbb{R}^D \rightarrow \mathbb{R}$$

Mixture of Gaussians as an EBM

$$p_\theta(x) = \frac{\exp\left(-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right)}{1 + \sum_{i=2}^K \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right)}$$

$$p_\theta(x) = \frac{1}{1 + \sum_{i=2}^K \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right)} \cdot \frac{1}{\sum_{i=1}^K \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right)}$$

$$= \frac{1}{1 + \sum_{i=2}^K \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right)} \cdot \exp\left(-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right) - \log \sum_{i=1}^K \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right)$$

7. Benefits

- Hybrid Model for Classification and Conditional Generation
- Robustness against adversarial attacks
- Better Calibration

