

Types of learning: $\begin{cases} \rightarrow \text{supervised} \\ \rightarrow \text{unsupervised} \end{cases} \rightarrow \text{reinforcement}$

learn $p_\theta(y|x)$ \rightarrow

Supervised learning: given $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$, learn $f_\theta(x) \approx y$

Predicting probabilities often makes more sense than predicting labels.

$p(x, y) = p(x)p(y|x)$ (chain rule) \rightarrow must be positive, sum to 1.

$$\text{softmax}(f_1(x), \dots, f_n(x)) = \frac{e^{f_k(x)}}{\sum_{i=1}^n e^{f_i(x)}} = p(y=k|x)$$

The ML Method (Example in lecture: Logistic Regression)

1. Define your model class (how to represent the "program")
2. Define your loss function (how to determine "better" model)
3. Pick your optimizer (how to find "best" model)
4. Run on big CPU

$$(x, y) \sim p(x, y) \quad p(\theta) = \prod_i p(x_i, y_i) \quad (\text{i.i.d.}) = \prod_i p(x_i) p(y_i | x_i)$$

ideal \rightarrow $\log p(\theta) = \sum_i \log p(x_i) + \log p_\theta(y_i | x_i) = \sum_i \log p_\theta(y_i | x_i) + \text{const}$

$\theta^* \leftarrow \arg \max_{\theta} \sum_i \log p_\theta(y_i | x_i)$ maximum likelihood est. (MLE)

$\theta^* \leftarrow \arg \min_{\theta} \left(-\sum_i \log p_\theta(y_i | x_i) \right)$ negative log-likelihood (NLL)
loss func.

cross-entropy = how similar are two distributions p_θ and p ? \sim NLL.

$$H(p, p_\theta) = -\sum_y p(y|x_i) \log p_\theta(y|x_i) \approx -\log p_\theta(y_i | x_i)$$

Optimization also: (one approach)

Gradient descent

1. Find direction v where $L(\theta) \downarrow$

1. Compute $\nabla_\theta L(\theta)$

2. $\theta \leftarrow \theta + \alpha v$ learning rate

2. $\theta \leftarrow \theta - \alpha \nabla_\theta L(\theta)$

$$\text{gradient: } \nabla_\theta L(\theta) = \begin{bmatrix} \frac{dL(\theta)}{d\theta_1} \\ \vdots \\ \frac{dL(\theta)}{d\theta_n} \end{bmatrix}$$

logistic Equation: $\frac{1}{1 + e^{-\theta^T x}}$ (aka sigmoid)

can examine over/under-fitting \rightarrow

Risk = $\mathbb{E}_{x \sim p(x), y \sim p(y|x)} [L(x, y, \theta)]$ (Theoretical)

Empirical Risk = $\frac{1}{n} \sum_{i=1}^n L(x_i, y_i, \theta) \approx$ Risk (Practical)