

Logistic Regression

Input: d -dimension feature vector $x \in \mathbb{R}^d$

Output: A class or label $l \in \{1, \dots, k\}$

Ex: Image classification:

Given input image $x \in \mathbb{R}^{n \times n}$, predict a label for the image (e.g. cat/dog)

MNIST: $\{0, \dots, 9\}$, CIFAR-10 $\{0, \dots, 9\}$

Ex: Binary Classification:

Given some medical data $x \in \mathbb{R}^n$

try to predict incidence of heart disease

label: $\{0/1\}$

Logistic Regression model

Given input feature vector $x \in \mathbb{R}^n$,
the model predicts a prob. distribution over
the labels $\{1, \dots, k\}$

$$\text{softmax}(\tilde{A}x + b)$$

$\tilde{A} \in \mathbb{R}^{k \times n}$ $b \in \mathbb{R}^n$ $\in \mathbb{R}^k$

$$\underbrace{\text{affine linear map}: \mathbb{R}^n \rightarrow \mathbb{R}^k}_{\text{soft max}} \xrightarrow{\text{soft max}} P(\{1, \dots, k\})$$

$$\text{softmax: Input: } y \in \mathbb{R}^k = \begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix} \quad e^{y_j}$$

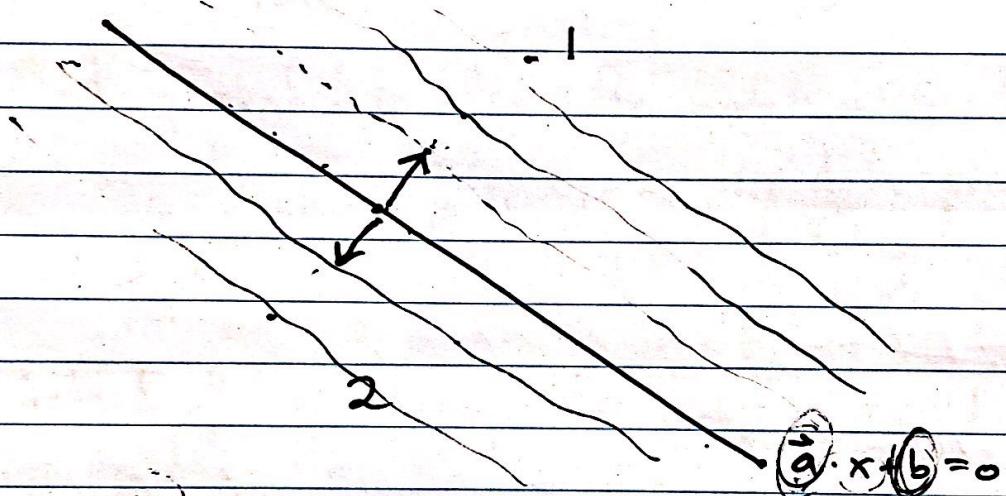
$$\text{Output: Distribution } P(j) = \frac{e^{y_j}}{\sum_{i=1}^k e^{y_i}}$$

"Take exponential and normalize"

Given input $x \in \mathbb{R}^n$, we return the distribution $\text{softmax}(Ax + b)$, A, b are parameters.

Ex: Two classes: Parameters are $\vec{a} \in \mathbb{R}^n$, $b \in \mathbb{R}$.

$$p(1) = \frac{e^{\vec{a} \cdot x + b}}{e^{\vec{a} \cdot x + b} + 1}, \quad p(2) = \frac{1}{e^{\vec{a} \cdot x + b} + 1}$$



$$\frac{p(1)}{p(2)} = e^{\vec{a} \cdot x + b}, \quad \log\left(\frac{p(1)}{p(2)}\right) \leftarrow \text{logarithm of the odds}$$

Assumption: $\log\left(\frac{p(1)}{p(2)}\right)$ is linear in the feature vector.

Learning the parameters \vec{a}, b from data

$$\text{Data: } \{(x_1, l_1), \dots, (x_n, l_n)\} = D$$

↑ ↑
Feature vectors labels

Given this data D , how can we estimate A, b ?

Maximum Likelihood

Given x_i , model gives softmax $(\vec{a} \cdot x_i + b)$

$$\frac{e^{\vec{a}_1 \cdot x_i + b_1}}{\sum_{i=1}^k e^{\vec{a}_i \cdot x_i + b_i}}$$

What probability does the model assign to ℓ_i ?

$$\frac{e^{\vec{a}_{\ell_i} \cdot x_i + b_{\ell_i}}}{\sum_{i=1}^k e^{\vec{a}_i \cdot x_i + b_i}} \quad \begin{array}{l} \text{(data points are} \\ \text{indep. so so} \\ \text{take product)} \end{array}$$

Instead of considering this prob, we consider its negative logarithm:

$$\log \left(\sum_{i=1}^k e^{\vec{a}_i \cdot x_i + b_i} \right) - (\vec{a}_{\ell_i} \cdot x_i + b_{\ell_i})$$

Logistic Regression loss:

$$L_D(A, b) = \sum_{(x, \ell) \in D} \left[\log \left(\sum_{i=1}^k e^{a_i \cdot x + b_i} \right) - (\vec{a}_{\ell} \cdot x + b_{\ell}) \right]$$

$$(A, b) = \min_{A, b} L_D(A, b)$$

Basic Statistical Learning Theory

- Goal: Estimate an unknown probability dist. D on a set X from samples (i.i.d)
 $x_1, \dots, x_n \in X$.

- Introduce a family of distributions p_θ for $\theta \in \Theta$ and try to choose θ to "match" the samples.

- Maximum Likelihood Estimate: Choose θ to maximize the probability of the samples.

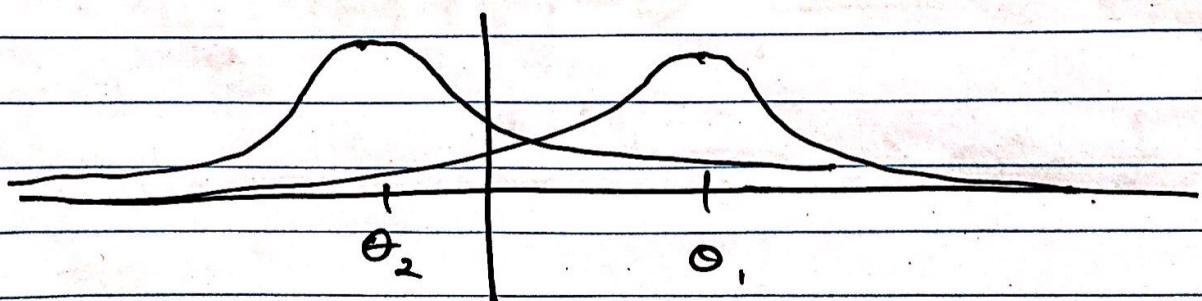
Ex: Let $X = \mathbb{R}$, have some samples

x_1, \dots, x_n drawn from a distribution D .

- Let p_θ be a Gaussian w/ variance 1, centered at $\theta \in \Theta = \mathbb{H}$, i.e.,

$$\text{density } p_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}}$$

centered at θ w/ variance 1



- Use the samples to find the center θ .

Maximum Likelihood Estimation

- Given $\Theta \in \mathbb{H} = \mathbb{R}$, what is the probability, of $\{x_j\}_{j=1}^n$?
Likelihood function (function of Θ)

- Samples independent: $p_\Theta(\{x_j\}_{j=1}^n) = \prod_{j=1}^n p_\Theta(x_j)$.

$$= \frac{1}{(2\pi)^{\frac{n}{2}}} \prod_{j=1}^n e^{-(x_j - \Theta)^2/2} = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\sum_{j=1}^n (x_j - \Theta)^2/2}.$$

- MLE: Choose Θ to maximize this!

- Often it's useful to consider $\log(p_\Theta(\{x_j\}_{j=1}^n))$.
log likelihood function

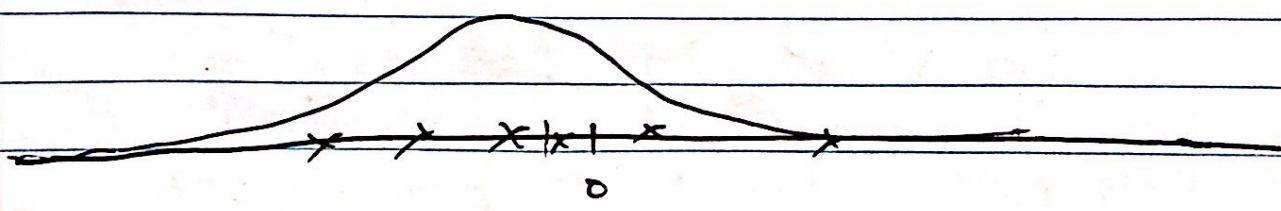
$$\Theta^* = \underset{\Theta \in \mathbb{H}}{\operatorname{argmax}} \log(p_\Theta(\{x_j\}_{j=1}^n)).$$

$$\left(\underset{\Theta \in \mathbb{H}}{\operatorname{argmin}} -\log(p_\Theta(\{x_j\}_{j=1}^n)) \right).$$

For this example: $\log(p_\Theta(\{x_j\}_{j=1}^n)) =$

$$= -\log(2\pi) \cdot \left(\frac{n}{2}\right) - \sum_{j=1}^n \frac{(x_j - \Theta)^2}{2}.$$

$$\Theta^* = \underset{\Theta \in \mathbb{R}}{\operatorname{argmin}} \sum_{j=1}^n \frac{(x_j - \Theta)^2}{2} \Rightarrow \Theta^* = \frac{1}{n} \sum_{j=1}^n x_j.$$



Classification / Logistic Regression

• For Classification $\tilde{X} = \underset{\text{features}}{X} \times \underset{\text{labels}}{Y}$

• We have samples $\{(x_j, y_j)\}_{j=1}^n$

• Suppose that D is an unknown distribution on X , but we're only trying to estimate $P_{\text{true}}(y|x)$ as a function of x .

• Introduce parameters $\Theta \in \mathbb{H}$, we define $p(y|x, \Theta)$ - called the model

• Now choose Θ to match the data $\{(x_j, y_j)\}_{j=1}^n$.

• MLE:

$$\Theta^* = \operatorname{argmax}_{\Theta \in \mathbb{H}} P(\{y_j\}_{j=1}^n | \{x_j\}_{j=1}^n, \Theta) =$$

$$= \operatorname{argmax}_{\Theta \in \mathbb{H}} \prod_{j=1}^n p(y_j|x_j, \Theta)$$

$$= \operatorname{argmin}_{\Theta \in \mathbb{H}} \sum_{j=1}^n -\log(p(y_j|x_j; \Theta)).$$

• Ex: Logistic Regression: $w \in \mathbb{R}^{k \times d}$ $b \in \mathbb{R}^d$.

- Model:

$$p(y|x, \Theta) = p(y|x, w, b)$$

- d -dimension of features, i.e. number of pixels
- k -number of classes.

$$\begin{aligned}
 p(y|x, w, b) &= \text{softmax} \left(\underbrace{\begin{pmatrix} w \in \mathbb{R}^{k \times d} \\ x \in \mathbb{R}^d \\ b \in \mathbb{R}^k \end{pmatrix}}_{\in \mathbb{R}^k} \cdot \underbrace{\begin{pmatrix} w_i x + b_i \\ \vdots \\ w_k x + b_k \end{pmatrix}}_{\in \mathbb{R}^k} \right) = \\
 &= \frac{\sum_{i=1}^k e^{w_i x + b_i}}{\sum_{i=1}^k e^{w_i x + b_i}} \cdot \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{\in \mathbb{R}^k} = \underbrace{\begin{pmatrix} e^{w_1 x + b_1} \\ \vdots \\ e^{w_k x + b_k} \end{pmatrix}}_{\in \mathbb{R}^k} \cdot \underbrace{\begin{pmatrix} y_i \\ \vdots \\ 0 \end{pmatrix}}_{\in \mathbb{R}^k} = \underbrace{\begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}}_{\in \mathbb{R}^k} \leftarrow \text{label } i
 \end{aligned}$$

• Need to calculate :

$$\begin{aligned}
 -\log(p(\hat{y}_i|x_i, w, b)) &= -\log \left(\frac{1}{\sum_{j=1}^k e^{w_j x_i + b_j}} e^{w_i x_i + b_i} \cdot y_i \right) \\
 &= \log(e^{w_i x_i + b_i} \cdot 1) - \log(e^{w_i x_i + b_i} \cdot y_i)
 \end{aligned}$$

$$(w, b)^* = \underset{\substack{w, b \\ \in \mathbb{R}^{k \times d} \times \mathbb{R}^k}}{\operatorname{argmin}} \sum_{i=1}^n \log(e^{w_i x_i + b_i} \cdot 1) - \log(e^{w_i x_i + b_i} \cdot y_i)$$

Bayesian Approach to Machine Learning

- Goal: Estimate an unknown distribution on X from data $\{x_j\}_{j=1}^n$.

- Build a model:

- Set of parameters (Θ)
- Family of distributions on X ,
 $p(x|\Theta)$ \leftarrow parameter $\in (\Theta)$

- Prior distribution on the parameters (Θ) .
 $q(\Theta)$.

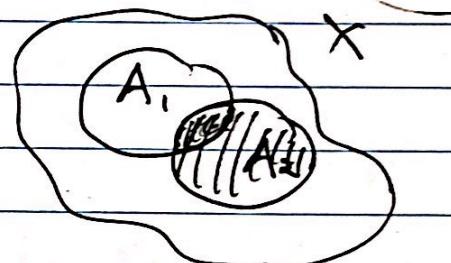
- Use Bayes' Law:

$$p(\Theta|x)p(x) = p(x \text{ and } \Theta) = p(x|\Theta)p(\Theta)$$

- Recall if A_1, A_2 events:

$$p(A_1|A_2) = \frac{p(A_1 \cap A_2)}{p(A_2)}$$

$$p(\Theta|x) = \frac{p(x|\Theta)q(\Theta)}{p(x)}$$



$$p(\Theta|x) \sim p(x|\Theta)q(\Theta)$$

\uparrow Posterior distribution \uparrow Likelihood function \uparrow Prior distribution

- Start with: prior $q(\Theta)$

- Add data x , replace q with the posterior

$$p(\Theta|x) \sim p(x|\Theta)q(\Theta)$$

- More data \rightarrow multiply by Likelihood & normalize.

- Left with a posterior distribution
 - Sample from posterior dist to approximate

$$P_{\text{pred}}(x) = \int_{\Theta} p(x|\Theta) p(\Theta) d\Theta$$

↑
posterior dist

- Choose Θ to maximize posterior distribution:

$$\Theta^* = \underset{\Theta \in \Theta}{\operatorname{argmax}} p(\Theta|x)$$

$$\Theta^* = \underset{\Theta \in \Theta}{\operatorname{argmin}} -\log(p(\Theta|x)) =$$

$$= \underset{\Theta \in \Theta}{\operatorname{argmin}} -\log(p(x|\Theta)) - \log(q(\Theta)) + \cancel{\log(p(\Theta))}$$

$$= \underset{\Theta \in \Theta}{\operatorname{argmin}} -\log(p(x|\Theta)) - \log(q(\Theta))$$

Negative log
likelihood

Regularization
coming from prior.

Ex: Image Classification / Logistic Regression

- Images $x \in X = \mathbb{R}^d$

- Labels $y \in Y = \{e_1, \dots, e_k\}$, $k = \# \text{ of labels}$

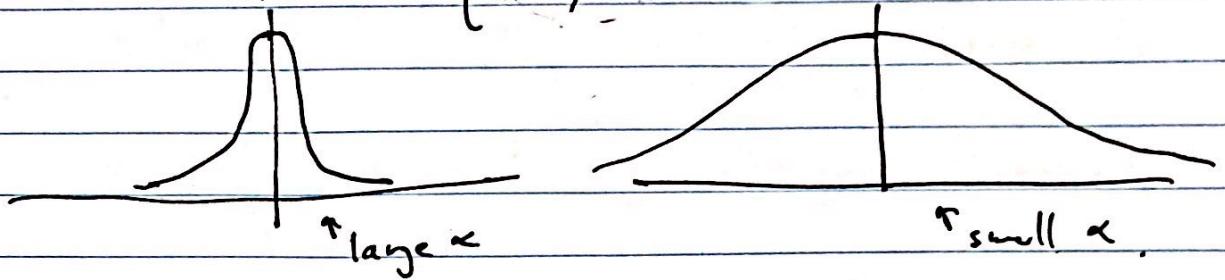
$$\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

- Data: $\{(x_j, y_j)\}_{j=1}^n$

- Model: $\Theta = (W, b) \in \mathbb{R}^{k \times d} \times \mathbb{R}^k$

$$p(y|x, \Theta) = \frac{e^{Wx+b} \cdot y}{e^{Wx+b} \cdot 1}$$

• Prior distribution $q(w, b) = C e^{-\alpha(\|w\|_2^2 + \|b\|_2^2)}$



- Calculate the posterior:

$$p(w, b | x, y) \propto \frac{p(y|x, w, b) p(w, b)}{p(y|x)}$$

$$\hookrightarrow = \int p(y|x, w, b) q(w, b) d\theta.$$

$$(w, b)^* = \underset{w, b}{\operatorname{argmin}} -\log(p(w, b | \{x_j, y_j\}_{j=1}^n)) =$$

$$= \underset{w, b}{\operatorname{argmin}} -\log(p(\{y_j\}_{j=1}^n | \{x_j\}_{j=1}^n, w, b)) - \log(q(w, b))$$

$$= \underset{w, b}{\operatorname{argmin}} -\log \left(\prod_{j=1}^n \underbrace{p(y_j | x_j, w, b)}_{e^{wx_j+b} \cdot y_j} \right) - \log \underbrace{(q(w, b))}_{C e^{-\alpha(\|w\|_2^2 + \|b\|_2^2)}}$$

$$= \underset{w, b}{\operatorname{argmin}} \sum_{j=1}^n \underbrace{\log(e^{wx_j+b} \cdot 1)}_{+ \alpha(\|w\|_2^2 + \|b\|_2^2)} - \log(e^{wx_j+b} \cdot y_j)$$