

ASSIGNMENT 1: [IFT6390]

LÉA RICARD & JOSEPH D. VIVIANO

1. PROBABILITIES

$P(D)$ = the probability of having disease.

$P(T)$ = the probability of having a positive test result.

- 1.5% of women in their 40s have breast cancer, therefore $P(D) = 0.015$.
- 87% true positive rate, therefore $P(T|D) = 0.87$.
- 9.6% false positive rate, therefore $P(T|\neg D) = 0.096$.

We want to know:

$$(1.1) \quad P(D|T) = \frac{P(T|D)P(D)}{P(T)}$$

Which means we need to calculate $P(T)$ which is $P(T|D) + P(T|\neg D) = 0.966$

$$(1.2) \quad P(D|T) = \frac{0.87 \times 0.015}{0.966} \approx 0.0135$$

2. CURSE OF DIMENSIONALITY

(A) Consider a hypercube in dimension d with side length c . What is the volume V ?

In the 2-dimensional case, $area = c^2$. In the 3 dimensional case, $V = c^3$. In the n -dimensional case, $V = c^d$.

(B) X is a random vector of dimension d ($x \in d$) distributed uniformly within the hypercube (the probability density $p(x) = 0$ for all x outside the cube). What is the probability density function $p(x)$ for x inside the cube? Indicate which property(ies) of probability densities functions allow you to calculate this result.

For all probability distributions:

$$(2.1) \quad \int_{-\inf}^{\inf} p(x) = 1$$

We know $p(x) = 0$ for all points outside of the hypercube. Therefore, $p(x) = 1$ for x inside the cube (since $1 - 0 = 1$).

(C) Consider the outer shell (border) of the hypercube of width 3% of c (covering the part of the hypercube extending from the faces of the cube and $0.03c$ inwards). For example, if $c = 100\text{cm}$, the border will be 3cm (left, right, top, etc ...) and will delimit this way a second (inner) hypercube of side $100 - 3 - 3 = 94\text{cm}$. If we generate a point x according to the previously defined probability distribution (by sampling), what is the probability that it falls in the border area? What is the probability that it falls in the smaller hypercube?

Let b be the amount to remove from the border on one side (i.e., left) of the outer hypercube.

$p(x_{\text{large}}) = 1$. Therefore in the general case the probability we are in the smaller hypercube is:

$$(2.2) \quad p(x_{\text{small}}) = (c - 2b)^d / c^d$$

$$(2.3) \quad p(x_{\text{border}}) = 1 - p(x_{\text{small}})$$

Therefore for the above example:

$$(2.4) \quad p(x_{\text{small}}) = (100 - 2 \times 3)^d / 100^d = 94^d / 100^d$$

And as before:

$$(2.5) \quad p(x_{\text{border}}) = 1 - p(x_{\text{small}})$$

And the probability we are in the border is:

(D) Calculate the above for $d = 1, 2, 3, 5, 10, 100, 1000$.

$$(2.6) \quad 1 - 94^1 / 100^1 = 0.06$$

$$(2.7) \quad 94^2 / 100^2 = 0.1163$$

$$(2.8) \quad 94^3 / 100^3 = 0.1694$$

$$(2.9) \quad 94^5 / 100^5 = 0.2661$$

$$(2.10) \quad 94^{10} / 100^{10} = 0.4614$$

$$(2.11) \quad 94^{100}/100^{100} = 0.9980$$

$$(2.12) \quad 94^{1000}/100^{1000} \approx 1(i.e., \text{tresgrand}).$$

(E) When the dimension grows, the probability that x falls into the narrow border at the edge of the hypercube becomes more likely, which is contrary of our intuitions at lower dimensions.

3. PARAMETRIC GAUSSIAN VS PARZEN WINDOW DENSITY ESTIMATION

3.1. Isotropic Gaussian Distribution. (A) The named parameters are μ , a d -long vector of means, and Σ , a $d \times d$ covariance matrix, where n is the number of data points.

(B)

$$(3.1) \quad \mu = \frac{1}{n} \sum_{i=1}^n x_i.$$

$$(3.2) \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T$$

(C) For the μ parameter, the algorithm complexity is in $\mathcal{O}(n)$, since it is summing over the vectors (x_i) .

For the Σ parameter, since the Gaussian density function is isotropic, the algorithm complexity is also in $\mathcal{O}(n)$. μ can be calculated by $\Sigma = \sigma I^2$.

(D)

$$(3.3) \quad p(x) = \frac{1}{(2\pi)^{\text{frac}{d}{2}} \sqrt{|\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

(E) The cost of this operation is $\mathcal{O}(1)$, since this is a straight calculation.

3.2. Parzen windows with Isotropic Gaussian Kernels. (A) The only thing learned during training of Parzen Window Density Estimation using Gaussian kernels is σ , since the peak of each Gaussian is each data point. Therefore, if σ is fixed by the user, nothing is learned during training (the Gaussians are simply centered on each training data point).

(B)

$$(3.4) \quad p(x) = \frac{1}{n} \sum_{i=1}^n \mathcal{N}_{x,\sigma}(x)$$

Expanded becomes:

$$(3.5) \quad p(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{(2\pi)^{\frac{d}{2}} \sigma^d} e^{-\frac{1}{2} \frac{d(x_{test}, x_{train})^2}{\sigma^2}}$$

(C) $\mathcal{O}(n)$, since we have to calculate the distance between x_{test} and all of the n data points x_{train} .

3.3. Capacity/Expressivity. (A) The Parzen Gaussian is more expressive, because it can store information for every data point. The capacity of the algorithm grows as we give it more data points, this isn't true for the Gaussian distribution, which averages over all data points, so it has a fixed capacity for a given dimensionality, no matter how many training data the algorithm is shown.

(B) Parzen windows with Isotropic Gaussian Kernels, in the case that we used a large number of training examples with a small σ would result in extreme memorization of the noise in the training data (i.e., overfitting).

(C) Because in parametric Gaussian density estimation, σ is learned from the data, while it is fixed for all data points when using Parzen windows.

4. SOFTMAX ACTIVATION FUNCTION

UNIVERSITÉ DE MONTRÉAL

Email address: joseph@viviano.ca