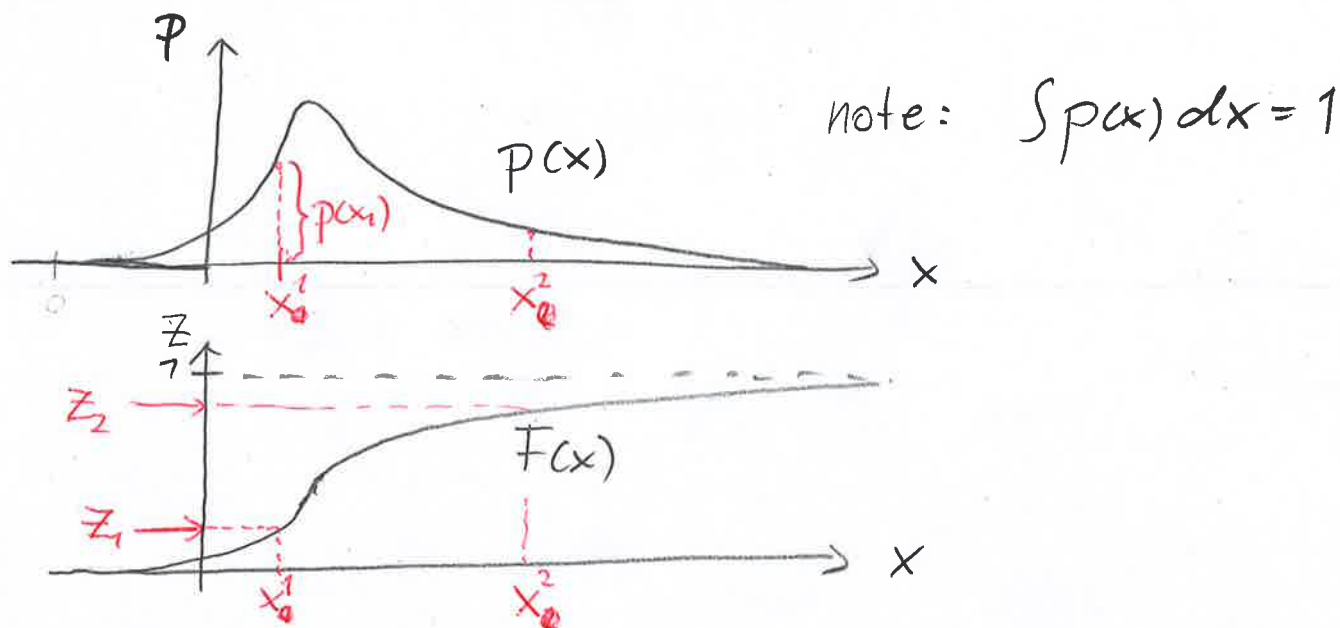


Blackboard 3.1: Random generation of data

You want to generate data points

$$\mathcal{X} = \{x^1, x^2, x^3, x^4\}$$

from a distribution $p(x)$



Q1: How do you do this on a computer?

- Integrate

$$\text{define } F(x) = \int_{-\infty}^x p(x') dx'$$

- Draw random numbers $z_k \in [0, 1]$
- $x^k = F^{-1}(z_k)$

Q2: What is the "likelihood" that you generate a point x^1 ?

$$P(x^1)$$

$$\boxed{\text{Note: Prob} = p(x) \cdot \underline{\Delta x}}$$

and all point x^1, x^2, x^3, x^4

$$P(\mathcal{X}) = P(x^1) \cdot P(x^2) \cdot P(x^3) \cdot P(x^4)$$

↑ ↑ ↑ ↑
independence

$(\Delta x)^4$

Blackboard 3.2: ML for Gaussian

A direct calculation

$$\begin{aligned} P_{\text{model}}(\mathcal{X} | x_{\text{center}}) &= P(x^1) \cdot P(x^2) \cdot \dots \cdot P(x^p) \\ &= \prod_k \left[\frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp \left\{ - \frac{(x^k - x_{\text{center}})^2}{2\sigma^2} \right\} \right] \\ &= \left[\frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \right]^k \exp \left\{ - \frac{1}{2\sigma^2} \sum_k (x^k - x_{\text{center}})^2 \right\} \end{aligned}$$

optimize parameter x_{center} :

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{\partial}{\partial x_{\text{center}}} P(\mathcal{X} | x_{\text{center}}) \\ &= \underbrace{P_{\text{model}}(\mathcal{X} | x_{\text{center}})}_{\neq 0} \cdot (-1) \frac{1}{2\sigma^2} \cdot 2 \cdot \sum_k (x^k - x_{\text{center}}) \\ \Rightarrow \quad \underline{\underline{x_{\text{center}}}} &= \frac{1}{p} \sum_{k=1}^p x^k \end{aligned}$$

B: alternatively with log-likelihood

$$\mathcal{L}(\mathcal{X} | x_{\text{center}}) = \ln P_{\text{model}}(\cdot | \cdot) = \sum_k \left[\underbrace{\ln \frac{1}{\sqrt{2\pi}} + \ln \frac{1}{\sigma}}_{\text{constant}} - \frac{(x^k - x_{\text{center}})^2}{2\sigma^2} \right]$$

Blackboard 3.3A : stochastic model
network output

$$(1a) \quad \hat{y}_{\vec{\omega}}(\vec{x}) = g^{(2)} \left[\sum_i \omega_i^{(2)} g^{(1)} \left(\sum_k \omega_{ik}^{(1)} x_k \right) \right]$$

generation of labels

$$(1b) \quad \hat{y}_{\vec{\omega}} = p(z=+1 | \vec{x})$$

↑ label generated by my model

A study point \vec{x}^n with $t^n = +1$
↑ label in data base

What is the probability that $(\vec{x}^n, +1)$
could have been generated by (1a), (1b)?

$$\begin{aligned} P(\vec{x}^n, z^n = +1) &= P(z^n = +1 | \vec{x}^n) \cdot P(\vec{x}^n) \\ &= \hat{y}_{\vec{\omega}}(\vec{x}^n) \cdot P(\vec{x}^n) \end{aligned}$$

B study point \vec{x}^n with $t^n = 0$
↑ label in data base

$$\begin{aligned} P(\vec{x}^n, z^n = 0) &= P(z^n = 0 | \vec{x}^n) \cdot P(\vec{x}^n) \\ &= (1 - \hat{y}_{\vec{\omega}}(\vec{x}^n)) \cdot P(\vec{x}^n) \end{aligned}$$

Blackboard 3.3 B (continued)

likelihood that set of all points

$$\mathcal{X} = \{ (\vec{x}^\mu, t^\mu) ; 1 \leq \mu \leq P \}$$

could have been generated by model

$$\begin{aligned}
 P(\mathcal{X}) &= \left[\prod_{\substack{\vec{x}^\mu \in \mathcal{C} \\ \Leftrightarrow t^\mu = 1}} (\hat{y}_{\vec{\omega}}(\vec{x}^\mu)) \right] \cdot \left[\prod_{\substack{\vec{x}^\mu \in \mathcal{C} \\ \Leftrightarrow t^\mu = 0}} (1 - \hat{y}_{\vec{\omega}}(\vec{x}^\mu)) \right] \cdot \left[\prod_{\mu} P(\vec{x}^\mu) \right] \\
 &= \prod_{\mu} \left[(\hat{y}_{\vec{\omega}}(\vec{x}^\mu))^{t^\mu} \cdot (1 - \hat{y}_{\vec{\omega}}(\vec{x}^\mu))^{(1-t^\mu)} \right] \cdot \left[\prod_{\mu} P(\vec{x}^\mu) \right]
 \end{aligned}$$

log-likelihood

$$E(\vec{\omega}) = -\ln P(\mathcal{X}) = -LL_{\vec{\omega}}$$

$$= -\sum_{\mu=1}^P \left[t^\mu \cdot \ln(\hat{y}_{\vec{\omega}}(\vec{x}^\mu)) + (1-t^\mu) \ln(1 - \hat{y}_{\vec{\omega}}(\vec{x}^\mu)) \right]$$

~~$-\sum_{\mu} \ln P(\vec{x}^\mu)$~~
 constant,
 does not depend
 on $\vec{\omega}$

$E(\vec{\omega})$: cross-entropy error function

↑ minimize with respect to parameters $\vec{\omega}$

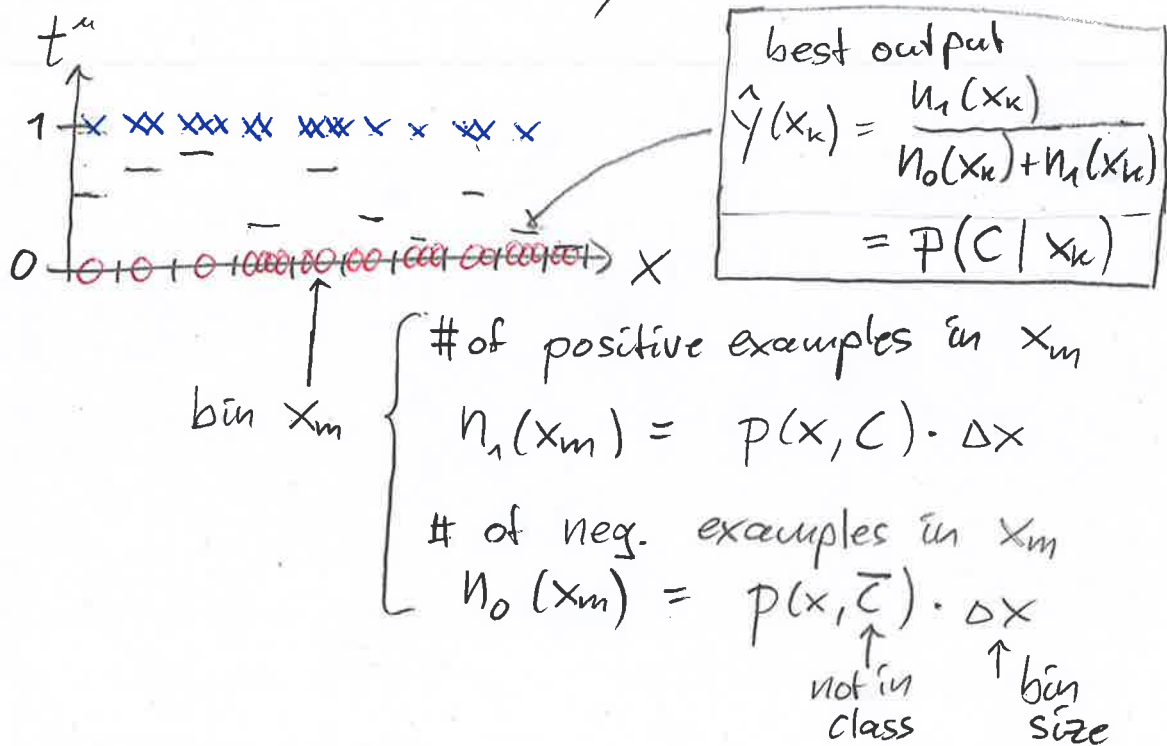
Blackboard 3.4 : output = probability?

start with

$$E = - \sum_{u=1}^P \left[t^u \cdot \ln \hat{y}^u + (1-t^u) \ln (1-\hat{y}^u) \right]$$

$$= - \sum_{\vec{x}^u \in C} \ln \hat{y}^u - \sum_{\vec{x}^u \notin C} \ln (1-\hat{y}^u)$$

hypothesis A : we have many examples



rewrite

$$E = - \sum_{\text{bins } m} \left[n_1(x_m) \cdot \ln \hat{y}(x_m) + n_0(x_m) \cdot \ln (1-\hat{y}(x_m)) \right]$$

hypothesis B : network is flexible enough

\Rightarrow in each bin, $\hat{y}(x_m)$ is arbitrary \Rightarrow optimize!

$$0 = \frac{\partial E}{\partial \hat{y}(x_m)} = \frac{n_1(x_m)}{\hat{y}(x_m)} - \frac{n_0(x_m)}{1-\hat{y}(x_m)}$$

$$\Rightarrow 0 = n_1(x_m) \cdot (1-\hat{y}(x_m)) - \hat{y}(x_m) \cdot n_0(x_m) \Rightarrow \hat{y}(x_m) = \frac{n_1(x_m)}{n_0(x_m) + n_1(x_m)}$$

Blackboard 3.5: Sigmoidal

output as probability

$$\hat{y}_1 = P(C_1 | x) \stackrel{\text{Bayes}}{=} \frac{P(\vec{x} | C_1) \cdot P(C_1)}{P(x)}$$

$$= \frac{P(x, C_1)}{P(x)} = \frac{P(x, C_1)}{P(x, C_1) + P(x, \bar{C}_1)}$$

$$= \frac{1}{1 + \frac{P(x, \bar{C}_1)}{P(x, C_1)}} = \frac{1}{1 + b}$$

$$= \frac{1}{1 + e^{-a}}$$

$0 < b < \infty$
positive parameter,
harder to treat
analytically

$$a = \ln \left[\frac{P(\vec{x}, C_1)}{P(\vec{x}, \bar{C}_1)} \right]$$

$a =$ "log-probability ratio"

$$-\infty < a < \infty$$

↑ unconstrained
parameter

Blackboard 3.6: mutually exclusive classes

example: 4 symbols A, B, C, D

with prob

symbol

1-hot-code

P_A

$$A = \{1, 0, 0, 0\}$$

P_B

$$B = \{0, 1, 0, 0\}$$

P_C

$$C = \{0, 0, 1, 0\}$$

P_D

$$D = \{0, 0, 0, 1\}$$

arbitrary

$$\vec{E} = \{t_1, t_2, t_3, t_4\} \text{ "1-hot-coding"}$$

probability to gen. arbitrary symbol \vec{E}

$$P_{\vec{E}} = P_A^{t_1} \cdot P_B^{t_2} \cdot P_C^{t_3} \cdot P_D^{t_4} = \prod_i [P_i]^{t_i} \text{ (check for symbol "C")}$$

total probability to generate M observed target vectors

P^{tot}

$$\prod_{m=1}^M \prod_i [P_i^m]^{t_i^m}$$

↑ all outputs
↑ all patterns

neg. log-likelihood

$$E = -LL = -\ln P^{\text{tot}} = -\sum_{m=1}^M \sum_i t_i^m \ln [P_i^m]$$

Probabilities $\sum_i P_i^m = 1$

⇒ describe P_i^m by softmax!

$$Y_i^m = \frac{e^{a_i}}{\sum_k e^{a_k}}$$