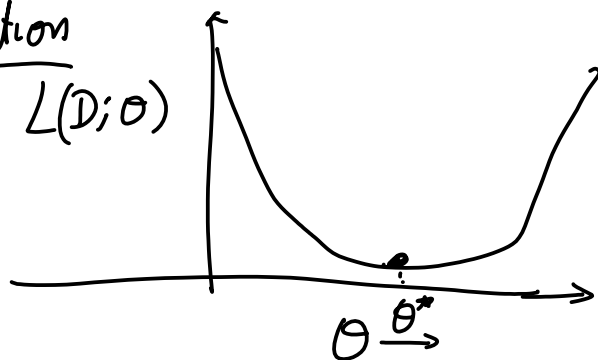


Machine learning as optimization

$$D = \{(\underline{x}_1, y_1) \dots (\underline{x}_n, y_n)\}$$

Training data



Choose a model

$\hat{y} = f(\underline{x}; \theta)$ = either a linear model
or a MLP

$l(\hat{y}, y)$ = error or loss for predicting \hat{y} when true value is y

$$L(D; \theta) = \sum_{i=1}^n l(\hat{y}_i, y_i) = \sum_{i=1}^n l(f(\underline{x}_i; \theta), y_i)$$

$$\theta^* = \arg \min_{\theta} \underbrace{L(D; \theta)}_{\text{Training loss}}$$

① Classify handwritten digits

↳ 6000 Handwritten images

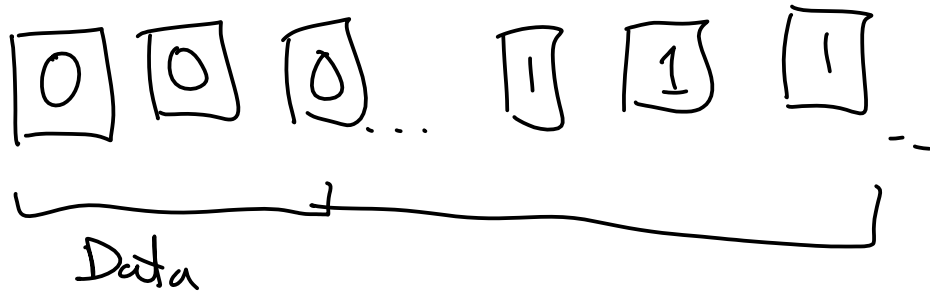
What matters is the loss/performance on previously unseen images

Data { Training data ✓ (only optimize on training data)
Test data (proxy for real world performance)

Probability is a measure/quantification of uncertainty

Sources of uncertainty in machine learning


eg
Handwritten
digits



we assume that all data is generated
by some data-generating process

↳ Incomplete observability

↳ Incomplete modeling

 \sim not taking factors into account


↳ Inherent uncertainty in the physical process

Probability theory

Def 1 Sample Space

Eg 2-coin tosses . H, T

$\{HH, HT, TH, TT\}$


Sample Space = Ω is the set of all possible outcomes

Eg Roll of a die with 6-sides

$$\Omega = \{1, 2, \dots, 6\}$$

Def 2 Event Space : is the space of potential results of the experiment

Eg Roll of dice

Event: whether we got dice ≥ 5

$$A = \{d \geq 5\} = \{5, 6\}$$

$$A \subseteq \Omega$$

Event Space is all possible values of such Events

Eg 2 coin toss $A = \{HH\} \subseteq \Omega$

$$A = \text{contains 1 tail} = \{HT, TT, TH\} \subseteq \Omega$$

For discrete sample space

Event space is the powerset of sample space

$$\text{Event Space} = \mathcal{A} = \mathcal{P}(\Omega) = 2^{\Omega}$$

↑ power set

$$\text{size of event space} = 2^{|\Omega|}$$

Def 3 Probability measure

$$P(A) \in [0, 1]$$

↑ Event

$$A \subseteq \Omega$$

$$A \in \mathcal{A}$$

↑ Event set

Def 4 (Sample space, Event Space, Probability measure)
 Probability space

Def 3 Probability measure is consistent if

Kolmogorov's axioms

- (a) $P(A) \in [0, 1]$
- (b) $P(\Omega) = 1$
- (c) $A_1 \cap A_2 = \emptyset$
 then $P(A_1) + P(A_2) = P(A_1 \cup A_2)$

countable infinite sets

$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n)$ when $A_n \cap A_{n2} = \emptyset$

e.g. dice
 $A_1 = \{5, 6\}$
 $A_2 = \{4, 3\}$
 $A_1 \cup A_2 = \{4, 3, 5, 6\}$

countable infinity — If you can assign a natural number to all

the elements

e.g. Natural numbers $\{1, \dots, \infty\}$

uncountable ∞ — Real numbers

Def 5 Random variable (RV)

e.g. 2 coin toss $\Omega = \{HH, HT, TH, TT\}$

RV maps an event $A \subseteq \Omega$ to a set of numbers $\left\{ \begin{array}{l} \text{Discrete} \\ \text{Continuous} \end{array} \right.$

$$X(\{HT, TH, TT\}) = \{1, 2, 3\}$$

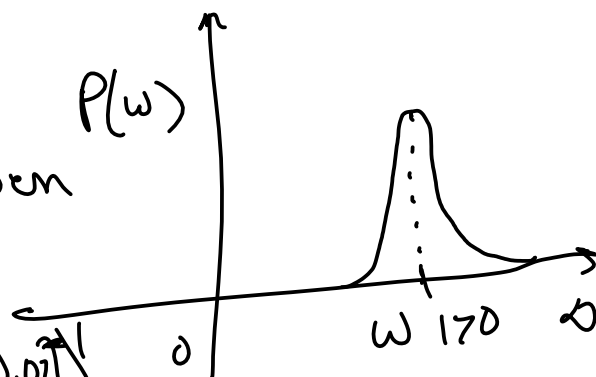
e.g. continuous RV : weight of mic

$$\Omega = [0, \infty)$$

closed interval

open

$$A \in [170, 171]$$



Event space of cont. RV $P(w \in [170.01, 170.02])$

\mathcal{A} = the set of all unions and intersections of countably infinite intervals (open, closed) } Borel set

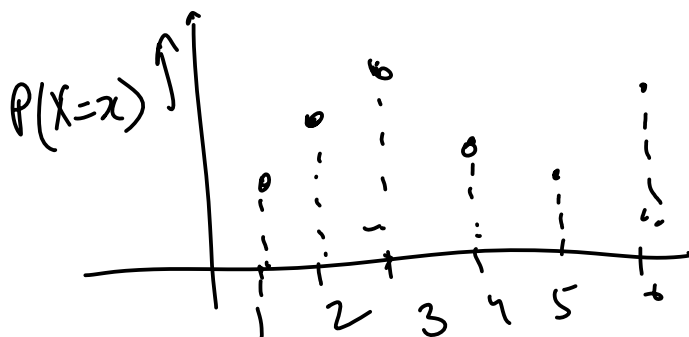
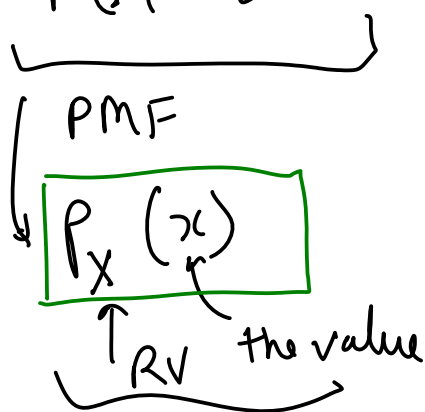
→ Different from the power set of all Real numbers

① Probability Mass function (PMF)
is only defined for discrete RV
 $\Omega = \{1, \dots, 6\}$

$$P(X=1) = p_1$$

$$P(X=2) = p_2$$

$$P(X=x) = p_x \quad \quad \quad \sum_{x \in \Omega} P(X=x) = 1$$



$P(x)$ Notational short cuts

Probability Density function (PDF)
only defined for continuous RV

$$P(X=x) = 0 \quad \text{if } x \in \mathbb{R} \quad X \text{ is continuous RV}$$

$$P(a \leq X \leq b) = P(X \in [a, b]) = \int_a^b \underbrace{f(x)}_{\text{PDF}} dx$$

A function $f: \mathbb{R} \rightarrow [0, \infty)$ is called PDF if

a) $f(x) \geq 0 \quad \forall x \in \mathbb{R}$

$$b) \int_{\mathbb{R}} f(x) dx = 1$$

$$c) P(a \leq X \leq b) = \int_a^b f(x) dx$$

$$\int_{170.1}^{170.2} 10 dx = 10(170.2 - 170.1) = 1$$

Common continuous distribution

$$N(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right)$$

$$x = \mu$$

$$\exp(0) = 1$$

$$= \frac{1}{\sqrt{2\pi} \sigma}$$

$$= 2$$

$$\sigma < \frac{1}{\sqrt{2\pi}}$$

$$\sigma = \frac{1}{2\sqrt{2\pi}}$$

Cumulative density function

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f(x) dx$$

Discrete

Continuous

$$F_X(x) = P(X \leq x) = \sum_{z \leq x} \underbrace{P(X=z)}_{\text{PMF}} \quad \left. \vphantom{\sum_{z \leq x}} \right\} \text{ For discrete RV}$$

Multi variable probability distribution

Joint probability mass function

$$P(X=x, Y=y) = P(X=x) \cap (Y=y)$$

e.g. 2 coin tosses = $X = \{H^0H^1, H^1T^2, T^3H^0, T^3T^3\}$
1 dice roll = Y

$$P(X=1, Y=5) = P(X=1 \text{ and } Y=5)$$

$$P(\underline{X} = \begin{bmatrix} x \\ y \end{bmatrix}) = P(\underbrace{X_1}_{\text{first element of Random vector}} = x \text{ and } X_2 = y)$$

Random variable
↓
Random vectors $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$
↑
RV

PDF = probability distribution function

$$P(a \leq X \leq b, c \leq Y \leq d) = \int_a^b \int_c^d \underbrace{f(x,y)}_{\text{PDF of 2 variables}} dx dy$$

$$\text{vectors } \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \underline{x}_{1:n/2} \\ \underline{x}_{n/2:n} \end{bmatrix}$$



$\{1, \dots, 6\}$

CDF

↳ Discrete $F_{xy}(x, y) = P(X \leq x, Y \leq y) = \sum_{a \leq x} \sum_{b \leq y} P(X = \overset{\downarrow}{a}, Y = \overset{\downarrow}{b})$

↳ Continuous $F_{xy}(x, y) = P(X \leq x, Y \leq y) = \int_{-\infty}^x \int_{-\infty}^y f(x, y) dy dx$