# 1 Physics Parallels

Quantum physics predictions are formulated in terms of probabilities. This principal limitation comes from existence of incompatible observables (measurable quantities) when it is not possible to measure the values of these quantities simultaneously.

Energy of the system plays crucial role in quantum and statistical physics.

- Lowest-energy state is the ground state

- Energy eigenstates (states with defined energy) are stationary states (measurement of observables is time independent)

- For some system, requirement of mathematical consistency, limits allowed values of energy eigenvalues i.e energy levels are discrete

Quantum physics generally deals with systems at subatomic level. Statistical physics derives its fundamental laws of macroscopic system, combining subatomic principles and general arguments translating extreme complexity of external interaction to the language of probability and statistics. In case of (quantum) statistical physics the inherent probabilistic nature of subatomic world is combined with probabilistic approach to, otherwise hopeless to solve, behavior of macroscopic bodies.

Important case is statistical equilibrium when the (macroscopic) system is in quantum stationary state. After many heuristics arguments (to be specified) we conclude

$$\log w_n = \alpha + \beta E_n \tag{1}$$

where $w_n$ is quantum probability of macroscopic system being in stationary state with energy $E_n$. The entropy can be expressed as average value of $\log w_n$ (to be explained)

$$S = -\sum_n w_n \log w_n \tag{2}$$

# 2 Definitions

Some definitions from Deep Learning Book using their LaTeXtemplates. Physics package provides the same symbols in more generic way eg. random vector $\boldsymbol{x}$ =\vb*{x} and vector $\mathbf{x}$ =\vb{x} instead of separate definition for vector $\boldsymbol{x}$ =\vx and random vector $\mathbf{x}$ =\rvx. Where possible I use symbols from `physics` package like $\|\boldsymbol{x}\|_p$ (\normp{x}{p}) for $L^p$
$L^p$ norm of $\boldsymbol{x}$

$$||\boldsymbol{x}||_p = \left(\sum_i |x_i|^p\right)^{\frac{1}{p}} \tag{3}$$

for $p \in \mathbb{R}, p \geq 1$. $L^2$ is abbreviated as $||\boldsymbol{x}|| = \|\boldsymbol{x}\|$. **Squared** $L^2$ norm meaning scalar product $\boldsymbol{x} \cdot \boldsymbol{x} = \boldsymbol{x}^T\boldsymbol{x} = \langle\boldsymbol{x}|\boldsymbol{x}\rangle$ (Dirac notation) doesn't have any special symbol.

$$L^\infty = \max_i |x_i| \tag{4}$$

# 3 MXNET functions

mxnet.gluon.loss.L2Loss : $L = \frac{1}{2}\sum_i |label_i - pred_i|^2$. While original definition was

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

# 4 Notation

## Numbers and Arrays

| | |
|---|---|
| $a$ | A scalar (integer or real) |
| $\boldsymbol{a}$ | A vector |
| $\boldsymbol{A}$ | A matrix |
| $\mathbf{A}$ | A tensor |
| $\boldsymbol{I}_n$ | Identity matrix with $n$ rows and $n$ columns |
| $\boldsymbol{I}$ | Identity matrix with dimensionality implied by context |
| $\boldsymbol{e}^{(i)}$ | Standard basis vector $[0, \ldots, 0, 1, 0, \ldots, 0]$ with a 1 at position $i$ |
| $\mathrm{diag}(\boldsymbol{a})$ | A square, diagonal matrix with diagonal entries given by $\boldsymbol{a}$ |
| $\mathrm{a}$ | A scalar random variable |
| $\mathbf{a}$ | A vector-valued random variable |
| $\mathbf{A}$ | A matrix-valued random variable |

## Sets and Graphs

| | |
|---|---|
| $\mathbb{A}$ | A set |
| $\mathbb{R}$ | The set of real numbers |
| $\{0, 1\}$ | The set containing 0 and 1 |
| $\{0, 1, \ldots, n\}$ | The set of all integers between 0 and $n$ |
| $[a, b]$ | The real interval including $a$ and $b$ |
| $(a, b]$ | The real interval excluding $a$ but including $b$ |
| $\mathbb{A} \backslash \mathbb{B}$ | Set subtraction, i.e., the set containing the elements of $\mathbb{A}$ that are not in $\mathbb{B}$ |
| $\mathcal{G}$ | A graph |
| $Pa_{\mathcal{G}}(\mathrm{x}_i)$ | The parents of $\mathrm{x}_i$ in $\mathcal{G}$ |

## Indexing

$a_i$ — Element $i$ of vector $\boldsymbol{a}$, with indexing starting at 1

$a_{-i}$ — All elements of vector $\boldsymbol{a}$ except for element $i$

$A_{i,j}$ — Element $i, j$ of matrix $\boldsymbol{A}$

$\boldsymbol{A}_{i,:}$ — Row $i$ of matrix $\boldsymbol{A}$

$\boldsymbol{A}_{:,i}$ — Column $i$ of matrix $\boldsymbol{A}$

$A_{i,j,k}$ — Element $(i, j, k)$ of a 3-D tensor $\mathsf{A}$

$\mathsf{A}_{:,:,i}$ — 2-D slice of a 3-D tensor

$\mathrm{a}_i$ — Element $i$ of the random vector $\mathbf{a}$

## Linear Algebra Operations

$\boldsymbol{A}^\top$ — Transpose of matrix $\boldsymbol{A}$

$\boldsymbol{A}^+$ — Moore-Penrose pseudoinverse of $\boldsymbol{A}$

$\boldsymbol{A} \odot \boldsymbol{B}$ — Element-wise (Hadamard) product of $\boldsymbol{A}$ and $\boldsymbol{B}$

$\det(\boldsymbol{A})$ — Determinant of $\boldsymbol{A}$

## Calculus

$\dfrac{dy}{dx}$ — Derivative of $y$ with respect to $x$

$\dfrac{\partial y}{\partial x}$ — Partial derivative of $y$ with respect to $x$

$\nabla_{\boldsymbol{x}} y$ — Gradient of $y$ with respect to $\boldsymbol{x}$

$\nabla_{\boldsymbol{X}} y$ — Matrix derivatives of $y$ with respect to $\boldsymbol{X}$

$\nabla_{\mathsf{X}} y$ — Tensor containing derivatives of $y$ with respect to $\mathsf{X}$

$\dfrac{\partial f}{\partial \boldsymbol{x}}$ — Jacobian matrix $\boldsymbol{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \to \mathbb{R}^m$

$\nabla_{\boldsymbol{x}}^2 f(\boldsymbol{x})$ or $\boldsymbol{H}(f)(\boldsymbol{x})$ — The Hessian matrix of $f$ at input point $\boldsymbol{x}$

$\displaystyle\int f(\boldsymbol{x})d\boldsymbol{x}$ — Definite integral over the entire domain of $\boldsymbol{x}$

$\displaystyle\int_{\mathbb{S}} f(\boldsymbol{x})d\boldsymbol{x}$ — Definite integral with respect to $\boldsymbol{x}$ over the set $\mathbb{S}$

## Probability and Information Theory

| | |
|---|---|
| a⊥b | The random variables a and b are independent |
| a⊥b \| c | They are conditionally independent given c |
| $P(a)$ | A probability distribution over a discrete variable |
| $p(a)$ | A probability distribution over a continuous variable, or over a variable whose type has not been specified |
| a $\sim P$ | Random variable a has distribution $P$ |
| $\mathbb{E}_{x \sim P}[f(x)]$ or $\mathbb{E}f(x)$ | Expectation of $f(x)$ with respect to $P(x)$ |
| $\text{Var}(f(x))$ | Variance of $f(x)$ under $P(x)$ |
| $\text{Cov}(f(x), g(x))$ | Covariance of $f(x)$ and $g(x)$ under $P(x)$ |
| $H(x)$ | Shannon entropy of the random variable x |
| $D_{\text{KL}}(P\|Q)$ | Kullback-Leibler divergence of P and Q |
| $\mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ | Gaussian distribution over $\boldsymbol{x}$ with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ |

## Functions

| | |
|---|---|
| $f : \mathbb{A} \to \mathbb{B}$ | The function $f$ with domain $\mathbb{A}$ and range $\mathbb{B}$ |
| $f \circ g$ | Composition of the functions $f$ and $g$ |
| $f(\boldsymbol{x}; \boldsymbol{\theta})$ | A function of $\boldsymbol{x}$ parametrized by $\boldsymbol{\theta}$. (Sometimes we write $f(\boldsymbol{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation) |
| $\log x$ | Natural logarithm of $x$ |
| $\sigma(x)$ | Logistic sigmoid, $\dfrac{1}{1 + \exp(-x)}$ |
| $\zeta(x)$ | Softplus, $\log(1 + \exp(x))$ |
| $\|\|\boldsymbol{x}\|\|_p$ | $L^p$ norm of $\boldsymbol{x}$ |
| $\|\|\boldsymbol{x}\|\|$ | $L^2$ norm of $\boldsymbol{x}$ |
| $x^+$ | Positive part of $x$, i.e., $\max(0, x)$ |
| $\mathbf{1}_{\text{condition}}$ | is 1 if the condition is true, 0 otherwise |

Sometimes we use a function $f$ whose argument is a scalar but apply it to a vector, matrix, or tensor: $f(\boldsymbol{x})$, $f(\boldsymbol{X})$, or $f(\mathbf{X})$. This denotes the application of $f$ to the array element-wise. For example, if $\mathbf{C} = \sigma(\mathbf{X})$, then $\mathsf{C}_{i,j,k} = \sigma(\mathsf{X}_{i,j,k})$ for all valid values of $i$, $j$ and $k$.

## Datasets and Distributions

| | |
|---|---|
| $p_{\text{data}}$ | The data generating distribution |
| $\hat{p}_{\text{data}}$ | The empirical distribution defined by the training set |
| $\mathbb{X}$ | A set of training examples |
| $\boldsymbol{x}^{(i)}$ | The $i$-th example (input) from a dataset |
| $y^{(i)}$ or $\boldsymbol{y}^{(i)}$ | The target associated with $\boldsymbol{x}^{(i)}$ for supervised learning |
| $\boldsymbol{X}$ | The $m \times n$ matrix with input example $\boldsymbol{x}^{(i)}$ in row $\boldsymbol{X}_{i,:}$ |