Part V: Probability and Statistics

- V.1 Mean, Variance, and Probability
- V.2 Probability Distributions
- V.3 Moments, Cumulants, and Inequalities of Statistics
- V.4 Covariance Matrices and Joint Probabilities
- V.5 Multivariate Gaussian and Weighted Least Squares
- V.6 Markov Chains

V.1 Mean, Variance, Probability

- Mean
 - Sample mean: done an experiment, got some output
 - Expected mean: known probabilities, but have not used them yet
 - Flip coins 0 or 1: sample mean → expected mean (approach with probability 1)
 - Law of large numbers: sample mean does approach ½ with probability 1 as the number of samples gets larger
- Variance
 - Sample variance: distance from the sample mean
 - Variance: expectation, not trial runs

- Mean is the average value or expected value
- Variance measures the average squared distance from the mean
- Probabilities of n different outcomes are positive numbers p_1, \ldots, p_n adding to 1
- Law of large numbers
 - with probability 1, the sample mean will converge to its expected value E[x] as the sample size N increases (e.g., flip coins)

sample mean: $m = \mu = \frac{1}{N} (x_1 + x_2 + \dots + x_N)$ where N: actual output \leftarrow what we got expected value: $m = E[x] = p_1 x_1 + p_2 x_2 + \dots + p_n x_n = \mathbf{p} \cdot \mathbf{x} \approx \int x p(x) dx \leftarrow$ what to expect n: number of possible different outputs with their probabilities

sample values: (five random freshmen ages) 18, 17, 18, 19, 17 → sample mean: 17.8 probabilities: (ages in a freshmen class) 17(20%), 18(50%), 19(30%) → expected age: 18.1

Applied Mathematics for Deep Learning

Highlights of Linear Algebra - 3

Variance (around the mean)

- Variance σ² measures expected distance (squared) from the expected mean E[x]
- Sample variance S² measures actual distance (squared) from the actual sample mean

in a set in the standard deviation - or C

$$= \text{Square root is the standard deviation o or S}$$
sample variance: $S^{2} = \frac{1}{N-1} \Big[(x_{1} - m)^{2} + \dots + (x_{N} - m)^{2} \Big]$
 $(N-1)$? one degree of freedom is already accounted for in the sample mean
$$\sum (x_{i} - m)^{2} = \sum (x_{i}^{2} - 2mx_{i} + m^{2}) = \sum x_{i}^{2} - 2m\sum x_{i} + Nm^{2} = \sum x_{i}^{2} - Nm^{2}$$
variance: $\sigma = E \Big[(x - m)^{2} \Big] = p_{1} (x_{1} - m)^{2} + \dots + p_{n} (x_{n} - m)^{2}$
 $E \Big[(x - m)^{2} \Big] = E \Big[x^{2} - 2mx + m^{2} \Big] = E \Big[x^{2} \Big] - 2mE \big[x \Big] + m^{2} = E \Big[x^{2} \Big] - \big(E \big[x \big] \big)^{2} = \sum p_{i} x_{i}^{2} - \Big(\sum p_{i} x_{i} \Big)^{2}$

$$\begin{cases} \text{sample ages: } m = 17.8, \ S^{2} = \frac{1}{5-1} \Big[(0.2)^{2} + (-0.8)^{2} + (0.2)^{2} + (1.2)^{2} + (-0.8)^{2} \Big] = 0.7 \\ \text{probabilities of ages: } m = 18.1, \ \sigma^{2} = (0.2)(-1.1)^{2} + (0.5)(-0.1)^{2} + (0.3)(0.9)^{2} = 0.49 \\ \text{Highlights of Linear Algebra - 4} \end{aligned}$$

Continuous Probability Distributions

- Age: year(17, 18, 19: n=3) → days(17 ≤ x ≤ 20: continuous range), probability distribution p(x)
- Uniform distribution
 - All ages between those numbers are "equally likely"
 - Chance F(x) that a random freshmen has age less than x



Figure V.1: F(x) is the cumulative distribution and its derivative p(x) = dF/dx is the **probability density function (pdf)**. For this **uniform distribution**, p(x) is constant between 17 and 20. The total area under the graph of p(x) is the total probability F = 1.

p(x)dx: probability of a sample falling in between x and dx

$$p(x)dx = F(x+dx) - F(x) \rightarrow p(x) = \frac{dF}{dx}, \text{ (probability of } a \le x \le b) = \int_{a}^{b} p(x)dx = F(b) - F(a)$$

$$m = E[x] = \int xp(x)dx \left[= \int_{x=17}^{20} x\left(\frac{1}{3}\right)dx = 18.5 \right]$$

$$\sigma^{2} = E\left[(x-m)^{2}\right] = \int p(x)(x-m)^{2} dx \left[= \int_{x=17}^{20} \frac{1}{3}(x-18.5)^{2} dx = \frac{3}{4} \right]$$
uniform distribution for $0 \le x \le a \rightarrow$

$$\begin{cases} \text{density: } p(x) = \frac{1}{a}, \text{ cumulative: } F(x) = \frac{x}{a} \\ \text{mean: } m = \frac{a}{2}, \text{ variance: } \sigma^{2} = \int_{0}^{a} \frac{1}{a} \left(x-\frac{a}{2}\right)^{2} dx = \frac{a^{2}}{12} \end{cases}$$

- Normal/Gaussian distribution: Bell-shaped Curve
- Central Limit Theorem
 - The average of N samples of "any" probability distribution approaches a normal distribution as N $\rightarrow \infty$



standard normal distribution: symmetric $x = 0 \rightarrow \underbrace{m = 0, \sigma^2 = 1}_{N(0,1)} \rightarrow p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

$$\begin{cases} \text{total probability: } \int_{-\infty}^{\infty} p(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1 \\ \text{mean: } m = E[x] = \int_{-\infty}^{\infty} xp(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} xe^{-x^2/2} dx = 0 \\ \text{variance: } \sigma^2 = E[x^2] = \int_{-\infty}^{\infty} p(x)(x-0)^2 dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-x^2/2} dx = 1 \\ N(0,\sigma) \rightarrow \begin{cases} F(\sigma) - F(-\sigma) \approx 2/3 \\ F(2\sigma) - F(-2\sigma) \approx 0.95 \end{cases}$$
$$N(0,1) \rightarrow \left[x \xrightarrow{\text{shift}} x - m \xrightarrow{\text{stretch}} \frac{x-m}{\sigma} \right] \rightarrow N(m,\sigma) = p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/2\sigma^2} \end{cases}$$

Applied Mathematics for Deep Learning

Highlights of Linear Algebra - 8

N coin flips and N $\rightarrow \infty$

$$N = 4: \left(\frac{1}{2} + \frac{1}{2}\right)^4 = \frac{1}{16} + \frac{4}{16} + \frac{6}{16} + \frac{4}{16} + \frac{1}{16} = 1$$

For any $N, m = \frac{N}{2}$ and $\sigma_N^2 = \frac{N}{4}$
subtracting m is "centering" or "detrending"
dividing by σ is "normalizing" or "standardizing" $\} \rightarrow X = \frac{x - m}{\sigma} = \frac{x - \frac{N}{2}}{\frac{\sqrt{N}}{2}} \rightarrow \begin{cases} (\text{mean of } X) = 0 \\ (\text{variance of } X) = 1 \end{cases}$
 $p(x) = 1$
 $p(x) = 1$
 $\frac{p(x) = 1}{\frac{1}{2}}$ $\frac{1}{16} = \frac{1}{\frac{2N}{M + 0}}$ $\frac{1}{N/2} = \frac{1}{N/2}$ $\frac{N}{N}$ $\frac{1}{N}$ $\frac{1}{N}$ $\frac{1}{N} = \frac{1}{N}$ $\frac{1}{N} = \frac{1}{N}$ $\frac{1}{N} = \frac{1}{N/2}$ $\frac{1}{N} = \frac{1}{N} = \frac{1}{N/2}$ $\frac{1}{N} = \frac{1}{N} =$

- Accepting uncertainty in the inputs(b) and estimating the variance in the outputs(x)
 - How to estimate the variance?
 - Often probability distributions p(x) are not known
 - Try different input \rightarrow compute the outputs \rightarrow take an average
 - Monte Carlo approximates an expected value E[x] by a sample average: error ~ O(1/ \sqrt{N}), slow improvement

samples
$$X_1$$
 to X_N $m = \frac{X_1 + \dots + X_N}{N}$ $S^2 = \frac{(X_1 - m)^2 + \dots + (X_N - m)^2}{N - 1}$
sum of outputs x_i
times probabilities p_i $m = \sum_{i=1}^n p_i x_i$ $\sigma^2 = \sum_{i=1}^n p_i (x_i - m)^2$
integral of outputs x
with probability density $m = \int xp(x) dx$ $\sigma^2 = \int (x - m)^2 p(x) dx$

V.2 Probability Distributions

- Binomial: tossing a coin n times
- Poisson: rare events
- Exponential: forgetting the past
- Gaussian=Normal: averages of many tries
- Log-normal: logarithm has normal distribution
- Chi-squared: distance squared in n dimensions
- Multivariable Gaussian: probabilities for a vector

V.3 Two Great Inequalities in Statistics

Markov's inequality (only) applies when all X_i≥0

$$\operatorname{Prob}\left[X \ge a\right] \le \frac{\overline{X}}{a} = \frac{mean}{a} = \frac{E[X]}{a} \leftarrow \overline{X} = E[X] \begin{cases} = \sum_{\text{all } s} X(s) \text{ times } (\operatorname{Probability of } X(s)) \\ \ge \sum_{X(s) \ge a} X(s) \text{ times } (\operatorname{Probability of } X(s)) \\ \ge \sum_{X(s) \ge a} a \text{ times } (\operatorname{Probability of } X(s)) \end{cases}$$
$$a = 3, E[X] = 1, X_i = i, \text{ show } \operatorname{Prob}\left[X \ge 3\right] \le \frac{1}{3}$$
$$0 p_0 + 1p_1 + 2p_2 + 3p_3 + 4p_4 + 5p_5 + \dots = 1 \rightarrow 0p_0 + 1p_1 + 2p_2 + 3(p_3 + p_4 + p_5 + \dots) + p_4 + 2p_5 + \dots = 1 \end{cases}$$

 $0p_{0} + 1p_{1} + 2p_{2} + 3p_{3} + 4p_{4} + 5p_{5} + \dots = 1 \rightarrow 0p_{0} + 1p_{1} + 2p_{2} + 3(p_{3} + p_{4} + p_{5} + \dots) + p_{4} + 2p_{5} + \dots = 1$ when $3(p_{3} + p_{4} + p_{5} + \dots) = 1$? $p_{1} = p_{2} = p_{4} = p_{5} = \dots = 0 \rightarrow p_{3} = 1/3, p_{0} = 2/3$ if there is equality and $Prob[X \ge a] = \frac{E[X]}{a}$, then all probabilities are actually zero except $E[X] = \frac{E[X]}{a}$

$$\operatorname{Prob}[X = a] = \frac{E[X]}{a} \text{ and } \operatorname{Prob}[X = 0] = 1 - \frac{E[X]}{a}$$

• Chebyshev inequality (no assumption $X_i \ge 0$)

apply Markov inequality to $Y_i = (X_i - m)^2$ with same probability p_i $\overline{Y} = \sum p_i Y_i = \sum p_i (X_i - m)^2 = \sigma^2 \rightarrow |X_i - m| \ge a \rightarrow (X_i - m)^2 \ge a^2$ $\Rightarrow \operatorname{Prob} \left[Y \ge a^2 \right] \le \frac{\overline{Y}}{a^2} = \frac{mean}{a^2} = \frac{\sigma^2}{a^2}$

V.4 Covariance Matrix

- Linear algebra: M different experiments at once •
 - Measure age, height, weight (a, h, w: M=3) of N people
 - M mean values and (separate) variances (\rightarrow matrix)
 - Connection between the M parallel experiments?
 - p_{ii} : probability that experiment 1 produces x_i and experiment 2 produces y_i COVa
- Example

ariance:
$$\sigma_{12} = \sum_{\text{all } i} \sum_{j} p_{ij} (x_i - m_1) (y_i - m_2)$$

- flip two coins separately vs. glue the coins together

for independent experiments, p_{ij} = probability of $(i, j) = (p_i)(p_j)$, V: covariance matrix

$$\begin{cases} \operatorname{coin} 1: \begin{pmatrix} x = 0 \text{ or } 1 \\ \operatorname{with} p = \frac{1}{2} \end{pmatrix} \rightarrow \begin{pmatrix} H & T \\ H & \frac{1}{4} & \frac{1}{4} \\ T & \frac{1}{4} & \frac{1}{4} \end{pmatrix} \text{ vs.} \begin{pmatrix} H & T \\ H & \frac{1}{2} & 0 \\ T & 0 & \frac{1}{2} \end{pmatrix} \rightarrow \begin{pmatrix} \sigma_{12} = 0 \\ V = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \end{pmatrix} \text{ vs.} \begin{pmatrix} \sigma_{12} = 1/4 \\ V = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \\ \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix} \\ \sigma_1^2 \sigma_2^2 = \sigma_1^2 \sigma_1^2 \sigma_2 \sigma_2^2 \end{bmatrix}$$

sample mean:
$$\overline{\mathbf{X}} = \frac{X_1 + \dots + X_N}{N} \rightarrow \text{sample covariance matrix: } \mathbf{S} = \frac{(X_1 - \overline{X})(X_1 - \overline{X})^T + \dots + (X_N - \overline{X})(X_N - \overline{X})^T}{N-1}$$

probability matrix : $\mathbf{P} = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix}$, total probability (all pairs) is 1: $\sum_{\text{all } i} \sum_j p_{ij} = 1$
row sum p_i of \mathbf{P} : $\sum_{j=1}^n p_{ij}$ = probability p_i of x_i in experiment 1
covariance matrix: $\mathbf{V} = \sum_{\text{all } i} \sum_j p_{ij} \begin{bmatrix} x_i - m_1 \\ y_i - m_2 \end{bmatrix} \begin{bmatrix} x_i - m_1 & y_i - m_2 \end{bmatrix} = \sum_{\text{all } i} \sum_j p_{ij} \begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_i - m_2) \\ (x_i - m_1)(y_i - m_2) & (y_i - m_2)^2 \end{bmatrix}$
 $= \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} = \mathbf{V}^T \rightarrow \text{at least positive semidefinite} \begin{bmatrix} p_{ij} (x_i - m_1)^2 \ge 0 \\ p_{ij} (y_i - m_2)^2 \ge 0 \rightarrow (\mathbf{V} \text{ is positive definite unless the} \\ \text{experiments are dependent} \\ \text{det}(\mathbf{V}_{ij}) = 0 \end{bmatrix}$
covariance matrix: $\sum_{M \times M} = E \Big[(\mathbf{X} - \overline{\mathbf{X}})(\mathbf{X} - \overline{\mathbf{X}})^T \Big] = \sum_j p_{ij} (\mathbf{X} - \overline{\mathbf{X}})(\mathbf{X} - \overline{\mathbf{X}})^T$

$$\begin{bmatrix} \text{mean of } \mathbf{c}^T \mathbf{X} : E[\mathbf{c}^T \mathbf{X}] = \mathbf{c}^T E[\mathbf{X}] = \mathbf{c}^T \overline{\mathbf{X}} \\ \text{variance of } \mathbf{c}^T \mathbf{X} : E[(\mathbf{c}^T \mathbf{X} - \mathbf{c}^T \overline{\mathbf{X}})(\mathbf{c}^T \mathbf{X} - \mathbf{c}^T \overline{\mathbf{X}})^T] = \mathbf{c}^T E[(\mathbf{X} - \overline{\mathbf{X}})(\mathbf{X} - \overline{\mathbf{X}})^T] \mathbf{c} = \mathbf{c}^T \mathbf{V} \mathbf{c} \ge 0 \\ \end{bmatrix}$$

 \rightarrow link between probability and linear algebra : $\mathbf{V} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ diagonalizing the covariance matrix \mathbf{V} means finding M independent experiments as combinations of the original M experiments

$$\begin{bmatrix} \text{mean and variance of } z = x + y \end{bmatrix}$$

$$\text{mean of sum = sum of mean : } \frac{1}{N} \sum_{i=1}^{N} (x_i + y_i) = \frac{1}{N} \sum_{i=1}^{N} x_i + \frac{1}{N} \sum_{i=1}^{N} y_i$$

$$E[x + y] = \sum_i \sum_j p_{ij} (x_i + y_j) = \sum_i \sum_j p_{ij} x_i + \sum_i \sum_j p_{ij} y_j = E[x] + E[y]$$

$$\sum_i \sum_j p_{ij} x_i = \sum_i (p_{i1} + \dots + p_{iN}) x_i = \sum_i p_i x_i = E[x]$$

$$\sigma_z^2 = \sum_i \sum_j p_{ij} (x_i + y_j - m_x - m_y)^2 = \sum_i \sum_j p_{ij} (x_i - m_x)^2 + \sum_i \sum_j p_{ij} (y_j - m_y)^2 + 2\sum_i \sum_j p_{ij} (x_i - m_x) (y_j - m_y)$$

$$= \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}$$

$$\mathbf{X} = \begin{bmatrix} x \\ y \end{bmatrix}, \ \mathbf{A} = \begin{bmatrix} 1 & 1 \end{bmatrix} \rightarrow \mathbf{Z} = \mathbf{A}\mathbf{X}$$

$$\sigma_z^2 = \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy} = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{A}\mathbf{V}\mathbf{A}^T$$

$$\mathbf{V}_{\mathbf{Z}} = \mathbf{A} \ \mathbf{V}_{\mathbf{X}} \ \mathbf{A}^T$$

$$_{M \times M}$$

[correlation ρ]

rescaling or standardizing the random variables *x* and *y* $\left(\text{like } \frac{v}{\|v\|} \right)$

$$\rightarrow X = \frac{x}{\sigma_x} \text{ and } Y = \frac{y}{\sigma_y} \text{ with } \sigma_X^2 = \sigma_Y^2 = 1$$
correlation of x and y $\left(\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}, -1 \le \rho_{xy} \le 1 \right)$ is the covariance of $X \left(= \frac{x}{\sigma_x} \right)$ and $Y \left(= \frac{y}{\sigma_y} \right)$

Applied Mathematics for Deep Learning

Highlights of Linear Algebra - 18