

Probability	
Experiment:	A repeatable procedure that generates an outcome
Sample space:	The set of all possible outcomes
Event:	An outcome or set of possible outcomes
Conditions for independence of A and B:	$P(A B) = P(A)$ $P(B A) = P(B)$ $P(A \cap B) = P(A)P(B)$
Bayes' Theorem:	$P(A B) = P(B A)P(A)/P(B)$
Mean	
Expectation	$E[X] = \sum x \cdot p(x)$
Variance	$Var(X) = E[(X - \mu)^2] = E[X^2] - \mu^2$
Standard deviation(s.d.)	$\sigma = \sqrt{Var(X)}$
E[aX + b]	$= a \cdot E[X] + b$
E[X + Y]	$= E[X] + E[Y]$
E[X - Y]	$= E[X] - E[Y]$
Variance	
Var(X)	$= E[(X - \mu)^2] = E[X^2] - \mu^2$
Var(nX)	$= E[(nX - n\mu)^2] = E[n^2(X - \mu)^2] = n^2Var(X)$
s.d.	$\sigma = \sqrt{Var(X)}$
Covariance and correlation	
Cov(X, Y)	$= E[XY] - \mu_X\mu_Y$
Cov(aX + b, cY + d)	$= ac \cdot Cov(X, Y)$
Cov(X + Y, Z)	$= Cov(X, Z) + Cov(Y, Z)$
Cov(X, X)	$= Var(X)$
For 2 dependent variables X and Y,	$Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$
For 2 independent variables X and Y,	$Cov(X, Y) = 0$
Cor(X, Y)	$= \rho = Cov(X, Y) / \sigma_X\sigma_Y$
Discrete r.v.	
Probability mass function (PMF):	$p_X(x) = P(X = x)$ for all values of x
Cumulative distribution function (CDF):	$F_X(x) = P(X \leq x)$ for all values of x
Continuous r.v.	
Probability density function (PDF):	$f(x)$
Cumulative distribution function (CDF):	$F(x)$
Distributions	
Bernoulli:	$X \sim Ber(p), E[X] = p, Var(X) = p(1 - p)$
Binomial:	$X \sim Bin(n, p), E[X] = np, Var(X) = np(1 - p)$
Geometric:	$X \sim Geo(p), E[X] = \frac{1}{p}, Var(X) = \frac{1-p}{p^2}$
Uniform:	$X \sim U(a, b), E[X] = \frac{a+b}{2}, Var(X) = \frac{1}{12}(b - a)^2$
Exponential:	$F(x) = 1 - e^{-\lambda x}, f(x) = \lambda e^{-\lambda x}$ for $x \geq 0$ $E[X] = 1/\lambda, Var(X) = 1/\lambda^2$
Normal:	$X \sim N(\mu, \sigma^2)$
Central Limit Theorem:	For any r.v. where n is large, sum (S_n) and average (\bar{X}_n) are approximately normal. $S_n \simeq N(n\mu, n\sigma^2)$ and $\bar{X}_n \simeq N(\mu, \frac{\sigma^2}{n})$
Linear Regression	
Least squares method:	Ensure sum of squared deviations is minimised $y = \beta_1 x + \beta_0$
Approximations:	$\widehat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}, \widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$
S_{xx}	$= \sum(x_i - \bar{x})^2$
S_{yy}	$= \sum(y_i - \bar{y})^2$
S_{xy}	$= \sum(x_i - \bar{x})(y_i - \bar{y})$
Sum of squared residuals (SSR)	$= \sum(y_i - \hat{y}_i)^2$, where \hat{y}_i is the approximated value of y_i

Goodness of fit $r^2 = 1 - \frac{SSR}{S_{yy}} = \frac{S_{xy}^2}{S_{xx}S_{yy}}$
Mean Absolute Error (MAE) $= \frac{1}{k} \sum \hat{y}_i - y_i $
Root Mean Square Error (RMSE) $= \sqrt{\frac{1}{k} \sum (\hat{y}_i - y_i)^2}$
Bayesian Inference
Hypothesis (H): A statement we wish to accept or reject
Prior $P(H)$: What we believe about the hypothesis without any evidence
Likelihood $L(H E) = P(E H)$: Likelihood of the hypothesis
Posterior $P(E H)$: Updated belief after seeing evidence
Conditional independence: $p(x y, \theta) = p(x \theta)$ $p(y x, \theta) = p(y \theta)$ $p(x, y \theta) = p(x \theta)p(y \theta)$, where θ is the parameter
Updating normal distributions: If n datapoints x_1, x_2, \dots, x_n are drawn from $N(\theta, \sigma^2)$, where θ^2 is known, $\sigma_{post}^2 = 1 / (\frac{1}{\sigma_{prior}^2} + \frac{n}{\sigma^2})$ $\mu_{post} = \frac{\sigma_{post}^2}{\sigma_{prior}^2} \mu_{prior} + \frac{n\sigma_{post}^2}{\sigma^2} \bar{x}$ Note also that $\frac{\sigma_{post}^2}{\sigma_{prior}^2} + \frac{n\sigma_{post}^2}{\sigma^2} = 1$, and $\sigma_{post}^2 < \sigma_{prior}^2$
Frequentist Inference
Type I error: Reject H_0 when H_0 is true
Type II error: Do not reject H_0 when H_1 is true
Significance level (α) $= P(\text{Type I error})$
Power of a test $= 1 - P(\text{Type II error})$ – note that this requires the knowledge of the distribution of our random variable under H_1 as well
Variance approximation: if σ^2 is unknown, we approximate it as $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ (i.e. $1/(n-1)$ * sample variance)
Conversion to T-distribution: Given $X \sim N(\mu_0, \sigma^2)$ with unknown σ^2 , $T = \frac{X - \mu_0}{\frac{\sqrt{s^2}}{\sqrt{n}}} \sim t(n-1)$
t-test for 2 samples: take pooled sample variance (s_p^2) as $s_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{(n-1) + (m-1)} \left(\frac{1}{n} + \frac{1}{m} \right)$
Maximum Likelihood Estimation
MLE: Consider data x drawn from some distribution with an unknown parameter p . The MLE estimation of p is the value of θ^* that maximises the likelihood $p(x p)$ $\theta^* = \arg \max_{\theta} P(x p = \theta)$
Geometric and Exponential: $\theta^* = \frac{\text{occurrences}}{\text{total time}} = \frac{1}{\text{average time per occurrence}}$
Bernoulli and Binomial: $\theta^* = \frac{\text{successes}}{\text{total}}$
Uniform: $\hat{b} = \max(x_1, \dots, x_n)$, $\hat{a} = \min(x_1, \dots, x_n)$
Normal: Sample mean and sample variance
Confidence Interval: Consider data drawn from some distribution with an unknown, fixed value θ . The interval estimator $[\hat{\theta}_L, \hat{\theta}_U]$ is called a confidence interval if $P(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$, where $1 - \alpha$ is the confidence level
Bias of point estimators: $\hat{\theta}$ is an unbiased estimator if $E[\hat{\theta}] = \theta$, and biased otherwise
Classification
Naïve Bayes': Assume that all features are mutually independent. $p(x_{i,1}, x_{i,2}, \dots, x_{i,m} y_i) = \prod p(x_{i,j} y_i)$
Laplacian smoothing: When the likelihood of any feature is 0, add 1 for each class and feature value
Prediction (1/0)

Ground truth (1/0)	True Positive (TP)	False Negative (FN)
	False Positive (FP)	True Negative (TN)
Accuracy: $\frac{TP + TN}{TP + TN + FP + FN}$		
Precision: $\frac{TP}{TP + FP}$		
Recall: $\frac{TP}{TP + FN}$		
F-score: Harmonic mean of precision and recall $F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$		