

# BAYESUVIUS

A VISUAL DICTIONARY OF BAYESIAN  
NETWORKS AND CAUSAL INFERENCE



ROBERT R. TUCCI

**Bayesuvius,**  
a visual dictionary of Bayesian Networks and  
Causal Inference

Robert R. Tucci  
[www.ar-tiste.xyz](http://www.ar-tiste.xyz)

January 22, 2024

This book is constantly being expanded and improved. To download the latest version, go to <https://github.com/rrtucci/Bayesuvius>

## **Bayesuvius**

by Robert R. Tucci

Copyright ©2020-2023, Robert R. Tucci.

This work is licensed under the Creative Commons Attribution-Noncommercial-No Derivative Works 3.0 United States License. To view a copy of this license, visit the link <https://creativecommons.org/licenses/by-nc-nd/3.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042.



Figure 1: View of Mount Vesuvius from Pompeii



Figure 2: Mount Vesuvius and Bay of Naples

# Contents

<b>Foreword</b>	<b>17</b>
<b>Appendices</b>	<b>18</b>
<b>A Navigating the ocean of Judea Pearl’s Books</b>	<b>19</b>
<b>B CI-2-3 track</b>	<b>20</b>
<b>C Notational Conventions and Preliminaries</b>	<b>24</b>
C.1 Some abbreviations frequently used throughout this book . . . . .	24
C.2 $\mathcal{N}(!a)$ . . . . .	24
C.3 Indicator function (a.k.a. Truth function) . . . . .	24
C.4 One hot vector . . . . .	25
C.5 $L^p$ norm . . . . .	25
C.6 Special sets . . . . .	27
C.7 Kronecker delta function . . . . .	27
C.8 Dirac delta function . . . . .	27
C.9 Majority function . . . . .	27
C.10 Underlined letters indicate random variables . . . . .	27
C.11 Probability distributions . . . . .	28
C.12 Discretization of continuous probability distributions . . . . .	28
C.13 Samples, i.i.d. variables . . . . .	29
C.14 Expected Value and Variance . . . . .	29
C.15 Conditional Expected Value . . . . .	30
C.16 Notation for covariances . . . . .	30
C.17 Conditional Covariance . . . . .	31
C.18 Normal Distribution . . . . .	32
C.19 Uniform Distribution . . . . .	33
C.20 Softmax function (a.k.a. Boltzmann Distribution) . . . . .	33
C.21 Sigmoid and log-odds functions . . . . .	34
C.22 Estimand, Estimator (curve-fit), Estimate, Bias . . . . .	35
C.23 Maximum Likelihood Estimate, Likelihood Ratio Test . . . . .	36
C.24 Mean Square Error (MSE) . . . . .	37
C.25 Cramer-Rao Bound . . . . .	39

C.26 Bayes Rule, Bayesian Updating And Conjugate Priors . . . . .	43
C.27 Linear regression, Ordinary Least Squares (OLS) . . . . .	44
C.27.1 LR, assuming $x_\sigma$ are non-random . . . . .	45
Derivation of LR From Minimization of Error . . . . .	46
Geometry of LR with non-random $x_\sigma$ . . . . .	47
LR Goodness of Fit, $R^2$ . . . . .	48
C.27.2 LR, assuming $x_\sigma$ are random . . . . .	50
Transforming expressions from non-random to random $x_\sigma$ . .	51
LR with random $x_\sigma$ , expressed in derivative notation . . . . .	53
Double regression of $y$ . . . . .	57
$R^2$ with random $x_\sigma$ . . . . .	59
C.28 Logistic Regression (LoR) . . . . .	60
C.29 Entropy, Kullback-Leibler divergence, Cross-Entropy . . . . .	60
C.30 Definition of various entropies used in Shannon Information Theory .	61
C.31 Mean log likelihood asymptotic behavior . . . . .	62
C.32 Arc Strength (Arc Force) . . . . .	64
C.33 Pearson Chi-Squared Test . . . . .	64
C.34 Demystifying Population and Sample Variances . . . . .	65
C.35 Independence of $\hat{\mu}$ and $\hat{\sigma}^2$ . . . . .	67
C.36 Chi-square distribution . . . . .	68
C.37 Student's t-distribution . . . . .	69
C.38 Hypothesis testing and 3 classic test statistics (Likelihood, Score, Wald)	72
C.39 Error Bars . . . . .	75
C.40 Confidence Interval . . . . .	76
C.41 Score p-value . . . . .	78
C.42 Convex/Concave functions, Jensen's Inequality . . . . .	80
C.43 Chebyshev's inequality . . . . .	81
C.44 Short Summary of Boolean Algebra . . . . .	83
C.45 Laplace transform . . . . .	84
C.45.1 Examples . . . . .	86
C.45.2 Properties . . . . .	87
C.46 Z-transform . . . . .	93
C.46.1 Examples . . . . .	96
C.46.2 Properties . . . . .	97
C.47 Legendre Transformation (dual functions) . . . . .	100
C.47.1 Examples . . . . .	101
C.47.2 Properties . . . . .	103
C.47.3 Connection to Fourier transform and Quantum Mechanics . .	106
C.48 Numpy tensor methods . . . . .	107

## D Definition of a Bayesian Network 113

<b>E</b>	<b>Bayesian Networks, Causality and the Passage of Time</b>	<b>117</b>
E.1	Unifying Principle of this book . . . . .	117
E.2	You say tomato, I say tomato . . . . .	118
E.3	A dataset is causal model free . . . . .	118
E.4	What is causality? . . . . .	119
E.5	Bayesian Networks and the passage of time . . . . .	120
E.6	Advice for the DAG-phobic . . . . .	121
<b>1</b>	<b>AdaBoost</b>	<b>122</b>
1.1	AdaBoost for general ensemble of w-classifiers . . . . .	122
1.2	AdaBoost for ensemble of tree stumps . . . . .	126
<b>2</b>	<b>ANOVA</b>	<b>128</b>
2.1	Law of Total Variance . . . . .	128
2.2	Sum of Squares Estimates . . . . .	129
2.3	F-statistic and hypothesis testing . . . . .	131
<b>3</b>	<b>ARACNE structure learning</b>	<b>133</b>
<b>4</b>	<b>Backdoor Adjustment Formula</b>	<b>135</b>
4.1	Examples . . . . .	136
<b>5</b>	<b>Back Propagation (Automatic Differentiation)</b>	<b>140</b>
5.1	Toy Example . . . . .	140
5.2	General Theory . . . . .	141
5.2.1	Jacobians . . . . .	141
5.2.2	Bnets for function composition, forward propagation and back propagation . . . . .	142
5.3	Application to Neural Networks . . . . .	144
5.3.1	Absorbing $b_i^\lambda$ into $w_{i j}$ . . . . .	144
5.3.2	Bnets for function composition, forward propagation and back propagation for NN . . . . .	145
5.4	General bnets instead of Markov chains induced by layered structure of NNs . . . . .	148
<b>6</b>	<b>Bell and Clauser-Horne Inequalities in Quantum Mechanics</b>	<b>149</b>
<b>7</b>	<b>Berkson's Paradox</b>	<b>150</b>
<b>8</b>	<b>Binary Decision Diagrams</b>	<b>152</b>
<b>9</b>	<b>Chow-Liu Trees and Tree Augmented Naive Bayes (TAN)</b>	<b>156</b>
9.1	Chow-Liu Trees . . . . .	156
9.2	Tree Augmented Naive Bayes (TAN) . . . . .	160

<b>10 Control Theory (linear, deterministic)</b>	<b>162</b>
10.1 Basic feedback model . . . . .	163
10.2 Classical model (analog) . . . . .	164
10.3 Modern model (analog) . . . . .	167
10.4 Classical model (digital) . . . . .	171
10.5 Modern model (digital) . . . . .	171
10.5.1 Discretizing derivatives . . . . .	172
10.5.2 Solving Difference Equation . . . . .	173
10.6 Higher than first order differential (or difference) equations . . . . .	175
10.6.1 Differential Equations . . . . .	175
10.6.2 Difference Equations . . . . .	176
10.7 Time-Invariance, Causality, Stability . . . . .	177
10.8 Controllability, Observability . . . . .	178
10.9 Signal Flow Graph . . . . .	178
<b>11 Copula</b>	<b>183</b>
11.1 Examples . . . . .	186
<b>12 Counterfactual Reasoning</b>	<b>189</b>
12.1 The 3 Rungs of Causal AI . . . . .	189
12.2 Do operator . . . . .	190
12.3 Imagine operator . . . . .	190
<b>13 Cross-Validation</b>	<b>194</b>
<b>14 DAG Extraction From Text (DEFT)</b>	<b>197</b>
<b>15 Dataset Shift and Batch Normalization</b>	<b>198</b>
15.1 Covariate Shift . . . . .	199
15.2 Concept Shift . . . . .	199
15.3 Batch Normalization . . . . .	200
<b>16 Decision Trees</b>	<b>201</b>
16.1 Transforming a dtree into a bnet . . . . .	203
16.2 Structure Learning for Dtrees . . . . .	204
16.2.1 Information Gain, Gini . . . . .	205
16.3 Information Gain Ratio . . . . .	208
16.3.1 Pseudo-code . . . . .	209
<b>17 Decisions Based on Rungs 2 and 3: COMING SOON</b>	<b>211</b>



<b>18 Difference-in-Differences</b>	<b>212</b>
18.1 John Snow, DID and a cholera transmission pathway . . . . .	212
18.2 PO analysis . . . . .	214
18.3 Linear Regression . . . . .	216
<b>19 Diffusion Models</b>	<b>219</b>
19.1 Bnet for DM . . . . .	219
19.2 Mean Values $M^{t-1}(x^t)$ and $M_{\theta}^{t-1}(x^t)$ . . . . .	222
19.3 Loss function $\mathcal{L}$ . . . . .	225
19.4 Algorithms for training and sampling DM . . . . .	227
<b>20 Digital Circuits</b>	<b>229</b>
20.1 Mapping any dcircuit to a bnet . . . . .	229
20.1.1 Option A of Fig.20.2 . . . . .	229
20.1.2 Option B of Fig.20.2 . . . . .	230
<b>21 Do Calculus</b>	<b>231</b>
21.1 3 Rules of Do Calculus . . . . .	234
21.2 Parent Adjustment Formula . . . . .	235
21.3 Backdoor Adjustment Formula . . . . .	237
21.4 Frontdoor Adjustment Formula . . . . .	238
21.5 Comparison of Backdoor and Frontdoor adjustment formulae . . . . .	239
21.6 Do operator for DEN diagrams . . . . .	240
<b>22 Do Calculus proofs</b>	<b>243</b>
<b>23 D-Separation</b>	<b>261</b>
<b>24 D-Separation in Quantum Mechanics</b>	<b>264</b>
<b>25 Dynamical Bayesian Networks</b>	<b>265</b>
<b>26 Expectation Maximization</b>	<b>267</b>
26.1 The EM algorithm: . . . . .	268
26.1.1 Motivation . . . . .	269
26.2 Minorize-Maximize (MM) algorithms . . . . .	269
26.3 Examples . . . . .	271
26.3.1 Gaussian mixture . . . . .	271
26.3.2 Blood Genotypes and Phenotypes . . . . .	272
26.3.3 Missing Data/Imputation . . . . .	274
<b>27 Factor Graphs</b>	<b>275</b>

<b>28 Frisch-Waugh-Lovell (FWL) theorem</b>	<b>278</b>
28.1 FWL, assuming $x^\sigma$ are non-random . . . . .	278
28.2 FWL, assuming $x^\sigma$ are random . . . . .	279
<b>29 Frontdoor Adjustment Formula</b>	<b>281</b>
29.1 Examples . . . . .	282
<b>30 G-formula (Sequential Backdoor Adjustment Formula)</b>	<b>283</b>
<b>31 Gaussian Nodes with Linear Dependence on Parents</b>	<b>287</b>
<b>32 Generalized Linear Model (GLM)</b>	<b>290</b>
32.1 Exponential Family of Distributions . . . . .	290
32.2 GLM . . . . .	292
<b>33 Generative Adversarial Networks (GANs)</b>	<b>297</b>
<b>34 Goodness of Causal Fit</b>	<b>302</b>
<b>35 Gradient Descent</b>	<b>303</b>
<b>36 Granger Causality</b>	<b>305</b>
<b>37 Hidden Markov Model</b>	<b>308</b>
37.1 Calculating $P(x_t, v^n)$ and $P(x_t, x_{t+1}, v^n)$ . . . . .	310
37.2 Calculating $\mathcal{F}_t$ and $\overline{\mathcal{F}}_t$ . . . . .	311
37.3 Calculating $P(x^n v^n)$ . . . . .	312
37.4 Calculating $P(v^n A, B, \pi)$ . . . . .	313
37.5 Calculating $\hat{x}^n$ (Viterbi algorithm) . . . . .	314
37.6 Calculating $\hat{A}, \hat{B}, \hat{\pi}$ (Baum-Welch algorithm) . . . . .	316
<b>38 Identification of do queries via LDEN diagrams</b>	<b>318</b>
<b>39 Influence Diagrams &amp; Utility Nodes</b>	<b>320</b>
<b>40 Instrumental Inequality and beyond</b>	<b>322</b>
40.1 I-inequality . . . . .	322
40.1.1 I-inequality for binary z,d,y . . . . .	324
40.2 Bounds on Effect of IV on treatment outcome y . . . . .	325
<b>41 Instrumental Variables</b>	<b>328</b>
41.1 $\delta$ with unmeasured confounder . . . . .	328
41.2 $\delta$ (with unmeasured confounder) can be inferred via IV . . . . .	329
41.3 More general bnets with IVs . . . . .	330
41.4 Instrumental Inequality . . . . .	331

<b>42 Jackknife Resampling</b>	<b>332</b>
42.1 Case $A = A^n(\vec{x}) = \frac{1}{n} \sum_{\sigma} x^{\sigma}$ . . . . .	334
<b>43 Junction Tree Algorithm</b>	<b>336</b>
<b>44 Kalman Filter</b>	<b>337</b>
44.1 Prediction Problem . . . . .	338
44.2 Solution . . . . .	339
44.3 Simple Example . . . . .	340
44.4 Invariants . . . . .	341
44.5 Derivation of Solution . . . . .	341
<b>45 LATE (Local Average Treatment Effect)</b>	<b>343</b>
<b>46 LDEN with feedback loops</b>	<b>349</b>
<b>47 Linear and Logistic Regression</b>	<b>356</b>
47.1 Generalization to $x$ with multiple components (features) . . . . .	358
47.2 Alternative $V(b, m)$ for logistic regression . . . . .	358
<b>48 Linear Deterministic Bnets with External Noise</b>	<b>360</b>
48.1 Example of LDEN diagram . . . . .	360
48.2 LDEN equations and their 2 solutions . . . . .	361
48.3 Fully connected LDEN diagrams . . . . .	362
48.3.1 Fully connected LDEN diagram with $nx = 2$ . . . . .	362
48.3.2 Fully connected LDEN diagram with $nx = 3$ . . . . .	364
48.3.3 Fully connected LDEN diagram with arbitrary $nx$ . . . . .	366
48.4 Not fully connected LDEN diagrams . . . . .	368
48.5 LDEN diagram with conditioned nodes . . . . .	369
48.6 SCuMpy . . . . .	369
48.7 Non-linear DEN diagrams . . . . .	369
<b>49 Marginalizer Nodes</b>	<b>371</b>
<b>50 Markov Blankets</b>	<b>373</b>
<b>51 Markov Chain Monte Carlo (MCMC)</b>	<b>375</b>
51.1 Inverse Cumulative Sampling . . . . .	375
51.2 Rejection Sampling . . . . .	377
51.3 Metropolis-Hastings Sampling . . . . .	378
51.4 Gibbs Sampling . . . . .	381
51.5 Importance Sampling . . . . .	382
<b>52 Markov Chains</b>	<b>384</b>

<b>53 Mediation Analysis</b>	<b>385</b>
<b>54 Mendelian Randomization</b>	<b>392</b>
<b>55 Message Passing and Bethe Free Energy</b>	<b>394</b>
55.1 2MRFs . . . . .	394
55.2 Message Passing Intuition . . . . .	395
55.3 $-\ln Z_\theta =$ Free Energy (FE) . . . . .	399
55.4 $-\ln Z_{\theta^*} =$ Minimum FE . . . . .	400
55.5 $-\ln Z_\theta^{tree} =$ Tree FE (a.k.a. Bethe FE) . . . . .	401
55.6 $-\ln Z_{\theta^*}^{tree} =$ Tree Minimum FE, and message passing . . . . .	402
<b>56 Message Passing, Pearl's theory</b>	<b>406</b>
56.1 Distributed Soldier Counting . . . . .	406
56.2 Spring Systems . . . . .	408
56.3 BP for Markov Chains . . . . .	408
56.4 BP Algorithm for Polytrees . . . . .	416
56.4.1 How BP algo for polytrees reduces to the BP algo for Markov chains . . . . .	419
56.5 Derivation of BP Algorithm for Polytrees . . . . .	420
56.6 Example of BP algo for a Tree . . . . .	423
56.7 Bipartite bnets . . . . .	427
56.8 BP for bipartite bnets (BP-BB) . . . . .	428
56.8.1 BP-BB and general BP agree on Markov chains . . . . .	430
56.8.2 BP-BB and general BP agree on tree bnets. . . . .	432
56.9 BP-BB and sum-product decomposition . . . . .	434
<b>57 Message Passing in Quantum Mechanics</b>	<b>435</b>
<b>58 Meta-learners for estimating ATE</b>	<b>436</b>
<b>59 Missing Data, Imputation</b>	<b>440</b>
59.1 Imputation via EM . . . . .	441
59.2 Imputation via MCMC . . . . .	444
59.3 Multiple Imputations . . . . .	445
<b>60 Modified Treatment Policy</b>	<b>446</b>
60.1 One time MTP . . . . .	446
60.2 $\Delta_{ c}$ estimand . . . . .	450
60.3 Estimates of $\Delta_{ c}$ . . . . .	453
60.3.1 Empirical estimate of $\Delta_{ c}$ . . . . .	453
60.3.2 OR estimate of $\Delta_{ c}$ . . . . .	453
60.4 Other Estimands besides $\Delta_{ c}$ . . . . .	456
60.5 Multi-time MTP . . . . .	456

<b>61 Monty Hall Problem</b>	<b>459</b>
<b>62 Multi-armed Bandits</b>	<b>461</b>
62.1 Bnet for MAB . . . . .	462
62.2 Reward functions . . . . .	464
62.3 Regret functions . . . . .	466
62.4 Strategies with random exploration . . . . .	467
62.4.1 $\epsilon$ -greedy algorithm . . . . .	467
62.4.2 $\epsilon_t$ -greedy algorithm . . . . .	468
62.5 Strategies with nonrandom exploration . . . . .	468
62.5.1 Upper Confidence Bounds (UCB) algorithms . . . . .	468
Frequentist UCB (UCB1) algorithm . . . . .	469
Bayesian UCB algorithm . . . . .	469
62.5.2 Thompson Sampling MAB (TS-MAB) algorithm . . . . .	471
Bnet for general TS-MAB algorithm . . . . .	471
TS-MAB algorithm with Beta agent and Bernoulli environment	472
TS-MAB algorithm, skeletal reprise . . . . .	473
62.5.3 Grad-MAB algorithm . . . . .	474
<b>63 Naive Bayes</b>	<b>477</b>
<b>64 Neural Networks</b>	<b>478</b>
64.1 Activation Functions $\mathcal{A}_i^\lambda : \mathbb{R} \rightarrow \mathbb{R}$ . . . . .	479
64.2 Weight optimization via supervised training and gradient descent . .	480
64.3 Non-dense layers . . . . .	482
64.4 Autoencoder NN . . . . .	484
<b>65 Noisy-OR gate</b>	<b>485</b>
65.1 3 ways to interpret the parameters $\pi_i$ . . . . .	486
<b>66 Non-negative Matrix Factorization</b>	<b>490</b>
66.1 Bnet interpretation . . . . .	490
66.2 Simplest recursive algorithm . . . . .	491
<b>67 Observationally Equivalent DAGs</b>	<b>492</b>
67.1 Examples . . . . .	492
<b>68 Omitted Variable Bias</b>	<b>495</b>
<b>69 Personalized Expected Utility</b>	<b>501</b>
69.1 Goal of PEU Theory . . . . .	502
69.2 Bnets for PEU Theory . . . . .	503
69.3 Bounds on $EU$ for unspecified bnet . . . . .	503
69.4 Bounds on $EU$ for specific bnet families . . . . .	506

<b>70 Personalized Treatment Effects</b>	<b>507</b>
70.1 Goal, Strategy and Rationale of PTE theory . . . . .	508
70.2 Bnets for PTE theory . . . . .	510
70.3 $ATE = PB - PH$ . . . . .	511
70.4 Probabilities Relevant to PTE theory . . . . .	512
70.5 Symmetry . . . . .	517
70.6 Linear Programming Problem . . . . .	518
70.7 Special constraints . . . . .	519
70.8 Matrix representation of probabilities . . . . .	522
70.9 Bounds on Exp. Probs. imposed by Obs. Probs. . . . .	525
70.10 Bounds on $PNS3$ for unspecified bnet . . . . .	527
70.11 Bounds on $PNS3$ for specific bnet families . . . . .	533
70.12 Bounds on $ATE$ imposed by Obs. Probs. . . . .	533
70.13 Bounds on $PNS$ in terms of $ATE$ and Obs. Probs. . . . .	533
70.14 Numerical Examples . . . . .	534
<b>71 Plate Notation</b>	<b>536</b>
<b>72 Potential Outcomes and Beyond</b>	<b>538</b>
72.1 $G$ and $G_{den}$ bnets, the starting point bnets . . . . .	539
72.2 $G$ bnet with nodes $y^\sigma(0), y^\sigma(1)$ added to it. . . . .	541
72.3 Expected Values of treatment outcome $y^\sigma$ . . . . .	543
72.4 Translation Dictionary . . . . .	543
72.5 $\mathcal{Y}_{ d,x} = \mathcal{Y}_{d x}$ (SUTVA) . . . . .	544
72.6 Conditional Independence Assumption (CIA) . . . . .	545
72.7 Treatment Effects . . . . .	545
72.8 Insights into what makes treatment effects equal and $\mathcal{Y}_{1 0} = \mathcal{Y}_1$ . . . .	548
72.9 $G_{do+}$ bnet . . . . .	549
72.10 $ACE = ATE$ . . . . .	550
72.11 Good, Bad Controls . . . . .	551
72.12 PO Confounder Sensitivity Analysis . . . . .	552
72.13 $(SDO, ATE)$ space . . . . .	554
72.14 Strata-Matching . . . . .	557
72.14.1 Exact strata-matching . . . . .	557
Estimates of Treatment Effects . . . . .	557
Example, estimation of treatment effects . . . . .	559
72.14.2 Approximate strata-matching . . . . .	561
72.14.3 Unbiased strata-matching estimates . . . . .	561
72.15 Propensities . . . . .	563
72.16 Propensity based estimates of treatment effects . . . . .	567
72.17 Positivity . . . . .	568
72.18 Multi-time PO bnets (Panel Data) . . . . .	569

<b>73 Program evaluation and review technique (PERT)</b>	<b>573</b>
73.1 Example . . . . .	575
<b>74 Random Forest and Bagging</b>	<b>579</b>
74.1 Bagging (with fully-featured bags) . . . . .	579
74.2 Bagging (with randomly-shortened bags) . . . . .	581
<b>75 Recurrent Neural Networks</b>	<b>582</b>
75.1 Language Sequence Modeling . . . . .	585
75.2 Other types of RNN . . . . .	585
75.2.1 Long Short Term Memory (LSTM) unit (1997) . . . . .	587
75.2.2 Gated Recurrence Unit (GRU) (2014) . . . . .	589
<b>76 Regression Discontinuity Design</b>	<b>591</b>
76.1 PO analysis . . . . .	591
76.2 Linear Regression . . . . .	593
<b>77 Regularization of Loss Functions</b>	<b>594</b>
77.1 $L^p$ norm ROLF . . . . .	595
77.1.1 $L^1$ norm ROLF can lead to sparsity . . . . .	595
77.1.2 $L^2$ norm ROLF for Least Squares . . . . .	597
77.2 Proximal functions . . . . .	598
77.3 Proximal ROLF . . . . .	600
77.4 Unobserved Nodes of a bnet . . . . .	601
<b>78 Reinforcement Learning (RL)</b>	<b>603</b>
78.1 Exact RL bnet . . . . .	606
78.2 Actor-Critic RL bnet . . . . .	608
78.3 Q function learning RL bnet . . . . .	610
<b>79 Reliability Box Diagrams and Fault Tree Diagrams</b>	<b>612</b>
79.1 Minimal Cut Sets . . . . .	618
<b>80 Restricted Boltzmann Machines</b>	<b>620</b>
<b>81 ROC curves</b>	<b>622</b>
81.1 Terminology Table Adapted from Wikipedia Ref.[146] . . . . .	625
<b>82 Scoring the Nodes of a Learned Bnet</b>	<b>627</b>
82.1 Probability Distributions and Special Functions . . . . .	628
82.2 Single node with no parents . . . . .	630
82.3 Multiple nodes with any number of parents . . . . .	632
82.4 Bayesian Scores . . . . .	634
82.5 Information Theoretic scores . . . . .	634

<b>83 Selection Bias Removal</b>	<b>636</b>
83.1 Pre and Post Switch Nodes . . . . .	637
83.2 Removing SB from passive query $P(y x)$ . . . . .	639
83.3 Removing SB from active query $P(y \mathcal{D}x)$ . . . . .	641
<b>84 Shannon Information Theory</b>	<b>643</b>
<b>85 Shapley Explainability</b>	<b>644</b>
85.0.1 Numerical examples of SHAP . . . . .	647
<b>86 Simpson’s Paradox</b>	<b>650</b>
86.1 Pearl Causality . . . . .	652
86.2 Numerical Example . . . . .	654
<b>87 Structure and Parameter Learning for Bnets</b>	<b>655</b>
87.1 Overview . . . . .	655
87.2 Score based SL algorithms . . . . .	657
87.3 Constraint based SL algorithms . . . . .	658
87.4 Pseudo-code for some bnet learning algorithms . . . . .	659
<b>88 Support Vector Machines And Kernel Method</b>	<b>661</b>
88.1 Learning Algorithm for SVM Classifier . . . . .	662
88.2 Linear (dot-product) Kernel . . . . .	663
88.3 Alternatives to Linear Kernel . . . . .	667
88.4 Random Forest and Kernel Method . . . . .	667
<b>89 Survival Analysis</b>	<b>668</b>
89.1 $S(t)$ estimates . . . . .	670
89.1.1 No-censoring estimate of $S(t)$ . . . . .	670
89.1.2 Kaplan-Meier estimate of $S(t)$ . . . . .	670
89.2 $\lambda(t)$ models . . . . .	675
89.2.1 $\lambda(t)$ independent of covariates $Z$ . . . . .	675
89.2.2 $\lambda(t)$ dependent on covariates $Z$ . . . . .	676
89.3 $S_0(t)$ estimates . . . . .	679
<b>90 Synthetic Controls</b>	<b>681</b>
90.1 PO analysis . . . . .	683
<b>91 Targeted Estimator</b>	<b>685</b>
91.1 Goal, Strategy, and Rationale of TE theory . . . . .	685
91.2 Functional Calculus . . . . .	687
91.3 Linear Approximation of $\Psi[P_N]$ . . . . .	689
91.4 ATE estimand . . . . .	690
91.5 ATE estimates . . . . .	691



91.5.1	$\Psi^E$	691
91.5.2	$\Psi^G$	692
91.5.3	$\Psi^{IPW}$	692
91.5.4	$\Psi^{LIPW}$	692
91.5.5	$\Psi^{LIPW++}$ (a.k.a. $\Psi^{TMLE}$ )	696
91.6	$\Psi^{TMLE}$ in practice	699
<b>92</b>	<b>Time Series Analysis: ARMA and VAR</b>	<b>701</b>
92.1	White noise	701
92.2	Backshift operator	702
92.3	Metrics	702
92.4	Definition of $ARMA(p, q)$ , $AR(p)$ and $MA(q)$ .	704
92.5	Solving $AR(p)$	706
92.6	Solving $MA(q)$	707
92.7	Solving $ARMA(p, q)$	708
92.8	Auto-correlation and partial auto-correlation	708
92.9	Generating function of auto-correlation	712
92.10	Impulse Response	713
92.11	$AR(p)$ and Yule-Walker equations	714
92.12	Forecasting	716
92.13	Model Learning	722
92.14	Differencing and $ARIMA(p, d, q)$	722
92.15	Parameter Learning	726
92.15.1	PL of $AR(p)$	726
92.15.2	PL of $MA(q)$	728
92.15.3	PL of $ARMA(p, q)$	731
92.16	$VAR(p)$	732
<b>93</b>	<b>Transfer Learning</b>	<b>733</b>
<b>94</b>	<b>Transformer Networks</b>	<b>735</b>
94.1	Recurrent Neural Net with Attention	736
94.1.1	Single Head Attention	736
94.1.2	Multi-Head Attention	739
94.2	Vanilla tranet	741
94.2.1	Single Head Attention	746
94.2.2	Multi-Head Attention	748
94.2.3	Encoder	750
94.2.4	Decoder	751
<b>95</b>	<b>Transportability of Causal Knowledge</b>	<b>755</b>

<b>96 Turbo Codes</b>	<b>759</b>
96.1 Decoding Algorithm . . . . .	762
96.2 Message Passing Interpretation of Decoding Algorithm . . . . .	764
<b>97 Uplift Modelling</b>	<b>765</b>
97.1 UP types . . . . .	765
97.2 Some Relevant Technical Formulas from Chapter 72 . . . . .	766
97.3 UP Analysis . . . . .	767
97.4 UP Decision Trees . . . . .	769
97.4.1 Appendix, connection between $\Delta_c$ and $\Delta_{clj}$ . . . . .	774
<b>98 Variational Bayesian Approximation for Medical Diagnosis</b>	<b>775</b>
<b>99 Variational Bayesian Approximation via <math>D_{KL}</math></b>	<b>779</b>
99.1 Free Energy $\mathcal{F}(\vec{x})$ . . . . .	781
<b>100 XGBoost</b>	<b>784</b>
100.1 Divergences . . . . .	784
100.2 Minimizing Cost function for single tree . . . . .	786
100.3 Leaf Splitting . . . . .	789
100.4 Pruning . . . . .	789
100.5 Feature Binning . . . . .	791
100.6 Final estimate of target attribute . . . . .	791
100.7 Bnet for XGBoost . . . . .	792
<b>101 Zero Information Transmission (Graphoid Axioms)</b>	<b>794</b>
101.1 Consequences of Eq.(101.5) . . . . .	795
<b>Bibliography</b>	<b>797</b>

# Appendices

# Chapter 94

## Transformer Networks

The primary reference for this chapter is Ref.[78]. Ref.[78] is the highly influential 2017 paper entitled “Attention is all you need” that introduced **Transformer Networks** (tranets) and Attention into the AI vernacular. Besides Ref.[78], I also read blog posts such as Ref.[28] and the Wikipedia article on tranet (Ref. [158]). For a complete list of the large number of excellent blog post that I read to learn this subject, see my open source software texnn (Ref.[77]).<sup>1</sup>

Transformer Networks (tranet) have been taking the fields of Natural Language Processing (NLP) and Large Language Models (LLM) by storm in recent years. They were introduced in 2017 and already are the basis of numerous LLMs. For example, BERT (Bidirectional Encoder Representations from Transformers) and ChatGPT (Generative Pre-trained Transformer), two tranet libraries that have been trained with huge databases such as all of English Wikipedia (2,500M words).

How well ChatGPT works was a huge surprise to most people, including experts in AI/ML. My conjecture is that this surprising LLM performance is due to causality. Let me explain. I believe tranets and the LLM that use them, are just curve-fitters (so are Least Squares, vanilla NNs, Convolutional NNs, etc.). But, we lucked out, because tranets are very good at fitting causal data, and the space of all human generated text, including math equations and computer code, is causality connected (i.e., has a **causally connected topology**).

Normally, tranets are drawn as box diagrams that are somewhat cryptic and ambiguous, at least to me. In this chapter, instead of drawing them as box diagrams, I represent them as causal DAGs (bnets). This makes their causal nature more explicit than the box diagrams, and, in my opinion, also makes them less ambiguous and more understandable than the box diagrams.

Recurrent Neural Nets (RNNs) are discussed in Chapter 75. tranets are quickly displacing RNNs, an older method, in NLP. tranets are better than RNNs for doing NLP in several important ways. Whereas RNNs analyze the tokens (words) of a

---

<sup>1</sup>texnn is Python software that I wrote specifically for drawing the bnets of this chapter, but later I generalized it to a stand-alone app that can draw any bnet (including SCMs, NNs and tranets), not just a tranet bnet.

sentence sequentially (like a Kalman Filter), tranets analyze them in parallel, and thus are more amenable to parallel computing. Also, because RNNs analyze the words of a sentence sequentially, they tend to give more importance to the end of a sentence than to its beginning. That's because RNNs start forgetting the beginning of a sentence by the time they reach its end, like a patient with Alzheimer's. tranets do not suffer from this malady.

Dynamical bnets are discussed in Chapter 25. In Chapter 75, we showed that RNNs are dynamical bnets. In this chapter we will show that tranets dynamical bnets too.

In this chapter, we will use the Numpy-like tensor notation discussed in Section C.48. In particular, note that  $[n] = [0 : n] = \{0, 1, \dots, n - 1\}$  and that  $T^{[n],[m]}$  is an  $n \times m$  matrix.

## 94.1 Recurrent Neural Net with Attention

### 94.1.1 Single Head Attention

Let

$\ell$  be the maximum number of words allowed in a sentence. Some words might be blanks (padding).

$d$  be the so called **hidden or embedding dimension**.

$e_\alpha^t \in \mathbb{R}^d$  be a  $d$  column vector for word  $\alpha \in [\ell]$  at time  $t$ .

$W_{\underline{q}}^t, W_{\underline{k}}^t, W_{\underline{v}}^t \in \mathbb{R}^{d \times d}$  be the weight matrices for time slice  $t$ . These matrices are learned by training the net. The letters  $Q, K, V$  stand for Query, Key and Value, respectively. They transform  $e_\alpha^t$  as follows

$$v_\alpha^t = W_{\underline{v}}^t e_\alpha^t \quad (94.1)$$

$$q_\alpha^t = W_{\underline{q}}^t e_\alpha^t \quad (94.2)$$

$$k_\alpha^t = W_{\underline{k}}^t e_\alpha^t \quad (94.3)$$

Fig.94.1 represents a tranet of a 3-word sentence as a dynamical bnet. The TPMs, printed in blue, for bnet Fig.94.1, are as follows:

$$P(v_\alpha^t | e_\alpha^t) = \mathbb{1}(v_\alpha^t = W_{\underline{v}}^t e_\alpha^t) \quad (94.4)$$

$$P(q_\alpha^t | e_\alpha^t) = \mathbb{1}(q_\alpha^t = W_{\underline{q}}^t e_\alpha^t) \quad (94.5)$$

$$P(k_\alpha^t | e_\alpha^t) = \mathbb{1}(k_\alpha^t = W_{\underline{k}}^t e_\alpha^t) \quad (94.6)$$

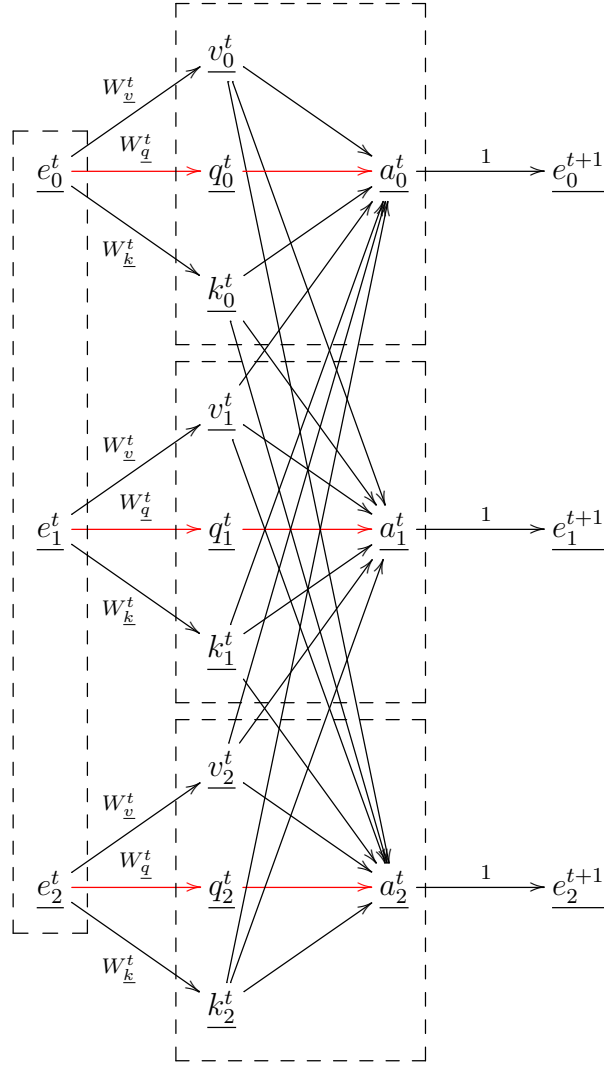


Figure 94.1: Dynamical bnet with single-head Attention for 3 words. Time-slice  $t$ . Note that  $k_\alpha^t$  for all  $\alpha$  points to  $\underline{a}_{\alpha'}^t$  for all  $\alpha'$ . Likewise,  $\underline{v}_\alpha^t$  for all  $\alpha$  points to  $\underline{a}_{\alpha'}^t$  for all  $\alpha'$ . However,  $\underline{q}_\alpha^t$  points only to  $\underline{a}_\alpha^t$ .

$$P(e_\alpha^{t+1} | a_\alpha^t) = \mathbb{1}(e_\alpha^{t+1} = a_\alpha^t) \quad (94.7)$$

$$P(a_\alpha^{t+1} | v_\alpha^t, q_\alpha^t, k_\alpha^t) = \mathbb{1}(a_\alpha^{t+1} = \sum_{\alpha' \in [\ell]} v_{\alpha'}^t P(\alpha' | \alpha)) \quad (94.8)$$

where the conditional probability  $P(\alpha' | \alpha)$  is called defined as

$$P(\alpha'|\alpha) = \text{softmax} \left[ \sum_{\delta \in [d]} (k^t)^{\delta, [\ell]} (q^t)^{\delta, \alpha} \right] (\alpha'|\alpha) \quad (94.9)$$

$$= \frac{e^{(k_{\alpha'}^t)^T q_{\alpha}^t}}{\sum_{\alpha'' \in [\ell]} e^{(k_{\alpha''}^t)^T q_{\alpha}^t}} \quad (94.10)$$

The right hand side of Eq.(94.8) constitutes an average over all the word vectors  $\{v_{\alpha}^t : \alpha \in [\ell]\}$  in a sentence. This average is called the **Attention** (for a single head).<sup>2</sup>

$$\text{Attention}^{\delta, \alpha} ((v^t)^{[d], [\ell]}, (k^t)^{[d], [\ell]}, (q^t)^{[d], [\ell]}) = \sum_{\alpha' \in [\ell]} (v^t)^{\delta, \alpha'} P(\alpha'|\alpha)$$

(94.11)

On first encounter, the structure of an Attention bnet seems a bit mysterious. Then one realizes that this is an old friend. If the dashed boxes in Fig.94.1 are each “shrunk” to single nodes, then it becomes a TAN Bayes Net. Each of the 3 subgraphs  $\underline{e}^t, (\underline{v}^t, \underline{q}^t, \underline{k}^t), \underline{a}^t$  also constitutes a TAN Bayes net.<sup>3,4</sup> In broad terms, Fig.94.1 can be described by saying that each word undergoes a special kind of 3-class (q,k,v) Naive Bayes classification, and the results of that classification are sent to the new version of each word (except the q class which only sends info to one word, not all of them).

It’s also useful to think of Attention as a filter with input signal  $(e^t)^{[d], [\ell]}$  and output signal  $(e^{t+1})^{[d], [\ell]}$ .

Fig.94.1 can be “folded” (i.e., the 3 words can be represented by as single node). When folded, Fig.94.1 becomes Fig.94.2. Note that in Fig.94.2, we have started indicating the shapes of tensors by a superscript, using the tensor notation explained in Section C.48. We will continue doing this henceforth in this chapter.

The structural equations for Fig.94.2, printed in blue, are as follows.

$$(a^t)^{[d], [\ell]} = \text{Attention}((v^t)^{[d], [\ell]}, (k^t)^{[d], [\ell]}, (q^t)^{[d], [\ell]}) \quad (94.12a)$$

$$(e^t)^{[d], [\ell]} = \text{prior} \quad (94.12b)$$

---

<sup>2</sup>Variations of this definition of Attention have been proposed. This particular one is the original one from the “Attention is all you need paper”. Some people call it the “scaled dot product Attention”.

<sup>3</sup>Tree Augmented Naive (TAN) Bayes nets were introduced in Chapter 9.

<sup>4</sup>A **reverse or upside down tree** is obtained by reversing the directions of all the arrows of a tree directed graph. A TAN Bayes net is normally defined as in Chapter9, as a Naive Bayes net augmented with a tree. In an Attention bnet, the Naive Bayes Net is augmented with a reverse tree (RT) instead of a tree (T), so technically Attention bnets contain RTAN Bayes nets, not TAN Bayes nets.

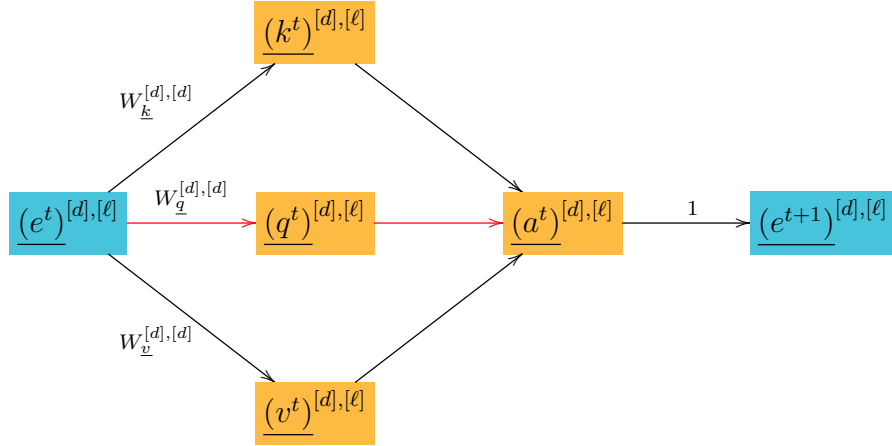


Figure 94.2: Folded version of Fig.94.1 when  $\ell = 3$ . Note that all orange nodes have the same tensor shape.

$$(e^{t+1})^{[d],[\ell]} = (a^t)^{[d],[\ell]} \quad (94.12c)$$

$$(k^t)^{[d],[\ell]} = W_{\underline{k}}^{[d],[d]}(e^t)^{[d],[\ell]} \quad (94.12d)$$

$$(q^t)^{[d],[\ell]} = W_{\underline{q}}^{[d],[d]}(e^t)^{[d],[\ell]} \quad (94.12e)$$

$$(v^t)^{[d],[\ell]} = W_{\underline{v}}^{[d],[d]}(e^t)^{[d],[\ell]} \quad (94.12f)$$

### 94.1.2 Multi-Head Attention

In this section, we will generalize the single head Attention, as defined in the previous section, to multi-head Attention.

Let

$n_{\underline{h}}$  = number of heads.  $\nu \in [n_{\underline{h}}]$ .

$d$  = same as before, the hidden, embedding dimension.  $\delta \in [d]$

$D = n_{\underline{h}}d$ .  $\Delta \in [D]$ . We will do some tensor reshaping:  $T^{[n_{\underline{h}}],[d]} \rightarrow T^{[D]}$ , or, in component form,  $T^{\nu,\delta} \rightarrow T^{\Delta}$ .

Consider weight matrices  $W_{\underline{k}}^{[D],[d]}$ ,  $W_{\underline{q}}^{[D],[d]}$ , and  $W_{\underline{v}}^{[D],[d]}$  such that



$$(k^t)^{\nu,\delta,\alpha} = \sum_{\delta' \in [d]} W_{\underline{k}}^{\nu,\delta,\delta'} (e^t)^{\delta',\alpha} \quad (94.13)$$

$$(q^t)^{\nu,\delta,\alpha} = \sum_{\delta' \in [d]} W_{\underline{q}}^{\nu,\delta,\delta'} (e^t)^{\delta',\alpha} \quad (94.14)$$

$$(v^t)^{\nu,\delta,\alpha} = \sum_{\delta' \in [d]} W_{\underline{v}}^{\nu,\delta,\delta'} (e^t)^{\delta',\alpha} \quad (94.15)$$

We define the **Multi-head Attention** by

$$\text{Attention}^{\nu,\delta,\alpha} ((v^t)^{[D],[\ell]}, (k^t)^{[D],[\ell]}, (q^t)^{[D],[\ell]}) = \sum_{\alpha' \in [\ell]} (v^t)^{\nu,\delta,\alpha'} P(\alpha'|\alpha, \nu) \quad (94.16)$$

where

$$P(\alpha'|\alpha, \nu) = \text{softmax} \left[ \sum_{\delta \in [d]} (k^t)^{\nu,\delta,[\ell]} (q^t)^{\nu,\delta,\alpha} \right] (\alpha'|\alpha, \nu) \quad (94.17)$$

$$= \frac{e^{(k^t)^{\nu,\delta,\alpha'} (q^t)^{\nu,\delta,\alpha}}}{\sum_{\delta \in [d]} \sum_{\alpha'' \in [\ell]} e^{(k^t)^{\nu,\delta,\alpha''} (q^t)^{\nu,\delta,\alpha}}} \quad (94.18)$$

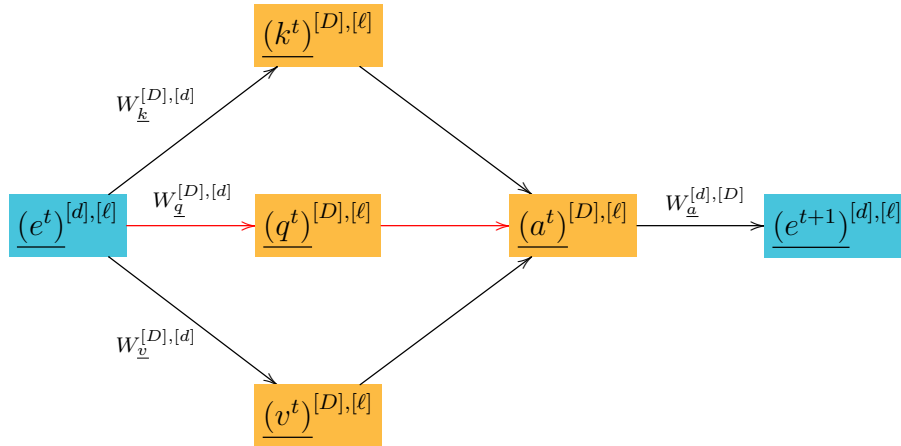


Figure 94.3: Dynamical bnet with single-head Attention for 3 words. Time-slice  $t$ . This is a generalization of the single head Attention of Fig.94.2. Note that all orange nodes have the same tensor shape.

The structural equations, printed in blue, for the bnet Fig.94.3, are as follows. Note that  $\text{Attention}()$  always has the same tensor shape as its 3 arguments. Note also

that the 3 weight matrices  $W_{\underline{k}}^{[D],[d]}$ ,  $W_{\underline{q}}^{[D],[d]}$ , and  $W_{\underline{v}}^{[D],[d]}$  raise the hidden dimension, whereas the weight matrix  $W_{\underline{a}}^{[d],[D]}$  lowers it.  $W_{\underline{a}}^{[d],[D]} = 1$  in the single head case.

$$(a^t)^{[D],[\ell]} = \text{Attention}((v^t)^{[D],[\ell]}, (k^t)^{[D],[\ell]}, (q^t)^{[D],[\ell]}) \quad (94.19a)$$

$$(e^t)^{[d],[\ell]} = \text{prior} \quad (94.19b)$$

$$(e^{t+1})^{[d],[\ell]} = W_{\underline{a}}^{[d],[D]}(a^t)^{[D],[\ell]} \quad (94.19c)$$

$$(k^t)^{[D],[\ell]} = W_{\underline{k}}^{[D],[d]}(e^t)^{[d],[\ell]} \quad (94.19d)$$

$$(q^t)^{[D],[\ell]} = W_{\underline{q}}^{[D],[d]}(e^t)^{[d],[\ell]} \quad (94.19e)$$

$$(v^t)^{[D],[\ell]} = W_{\underline{v}}^{[D],[d]}(e^t)^{[d],[\ell]} \quad (94.19f)$$

## 94.2 Vanilla tranet

In this section, we will discuss the tranet of the "Attention is all you need" paper, Ref.[78]. As is common in the literature, we will refer to that tranet as the "Vanilla" tranet. Ref.[78] describes its tranet graphically with Fig.94.4. Our goal is to find a causal DAG (bnet) version of that figure.

Let

$\ell$  = maximum number of words in a sentence segment.  $\alpha \in [\ell]$ ,  $\ell \sim 100$

$L$  = number of words in vocabulary,  $\beta \in [L]$ ,  $L \gg \ell$

$d = d_{\underline{q}} = d_{\underline{k}} = d_{\underline{v}} = 64$ , hidden dimension per head,  $\delta \in [d]$ .

$n_{\underline{h}} = 8$ , number of heads,  $\nu \in [n_{\underline{h}}]$

$D = n_{\underline{h}}d = 8(64) = 512$ , hidden dimension for all heads,  $\Delta \in [D]$

$\Lambda = 6$ , number of layers in plate (a.k.a., stack),  $\lambda \in [\Lambda]$

Our tensor notation is discussed in Section C.48. Here is a quick review of some of the more essential facts in that section on tensors:

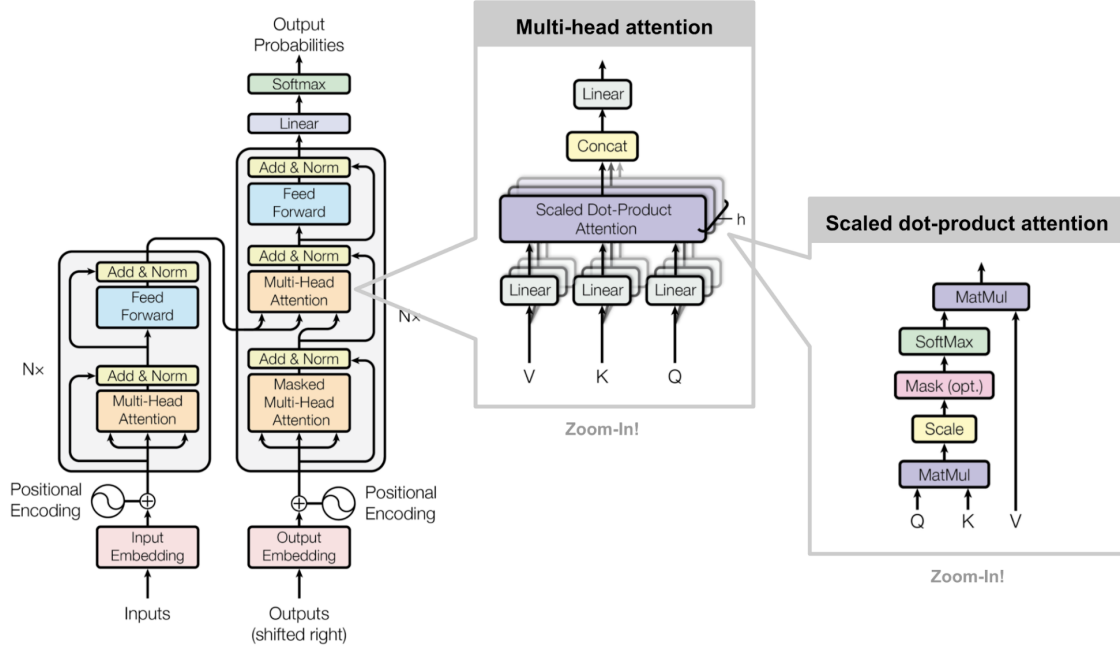


Figure 94.4: Vanilla tranet

- **reshaping**

$$T^{\nu, \delta} \rightarrow T^{\Delta} \quad (T^{[n_h], [d]} \rightarrow T^{[D]}) \quad (94.20)$$

$$T^{\Delta} \rightarrow T^{\nu, \delta} \quad (T^{[D]} \rightarrow T^{[n_h], [d]}) \quad (94.21)$$

- **concatenation**

$$T^{[n]} = (T^0, T^1, \dots, T^{n-1}) = (T^{\nu})_{\nu \in [n]} \quad (94.22)$$

- **Hadamard product (element-wise, entry-wise multiplication)**

$$T^{[n]} * S^{[n]} = (T^{\nu} S^{\nu})_{\nu \in [n]} \quad (94.23)$$

- **Matrix multiplication**

$T^{[n]} = T^{[n], [1]}$  is a column vector.

$$(T^{[n]})^T S^{[n]} = \text{scalar} \quad (94.24)$$

$$T^{[a], [b]} S^{[b], [c]} = \left[ \sum_{\beta \in [b]} T^{\alpha, \beta} S^{\beta, \gamma} \right]_{\alpha \in [a], \gamma \in [c]} \quad (94.25)$$

Most treatments of tranet, including the “Transformers is All You Need paper”, order the operations chronologically from left to right. So if  $A$  occurs before  $B$ , they write  $AB$ . This is contrary to what is done in Linear Algebra, where one orders the operations chronologically from right to left, and one writes  $BA$ . We will adhere to the Linear Algebra convention, since it is so prevalent and is the overwhelming precedent.

Before we present the bnet version of Fig.94.4, we discuss some of the definitions needed to understand and motivate Fig.94.4. Below will often accompany an equation in tensor component notation with, in parenthesis, the equivalent matrix equation.

- **Encoder Input**

$$x^{\beta,\alpha} = \delta(\beta, \beta(\alpha)) \left( x^{[L],[\ell]} \text{ has one hot columns.} \right) \quad (94.26)$$

- **Embedding (a.k.a. encoding) Matrix**

$$e^{\delta,\alpha} = \sum_{\beta} E^{\delta,\beta} x^{\beta,\alpha} \quad (e^{[d],[\ell]} = E^{[d],[L]} x^{[L],[\ell]}) \quad (94.27)$$

- **Weight matrices**

$$Q^{\nu,\delta,\alpha} = \sum_{\delta'} W_{\underline{q}}^{\nu,\delta,\delta'} e^{\delta',\alpha} \quad (Q^{[D],[\ell]} = W_{\underline{q}}^{[D],[d]} E^{[d],[\ell]}) \quad (94.28)$$

$$K^{\nu,\delta,\alpha} = \sum_{\delta'} W_{\underline{k}}^{\nu,\delta,\delta'} e^{\delta',\alpha} \quad (K^{[D],[\ell]} = W_{\underline{k}}^{[D],[d]} E^{[d],[\ell]}) \quad (94.29)$$

$$V^{\nu,\delta,\alpha} = \sum_{\delta'} W_{\underline{v}}^{\nu,\delta,\delta'} e^{\delta',\alpha} \quad (V^{[D],[\ell]} = W_{\underline{v}}^{[D],[d]} E^{[d],[\ell]}) \quad (94.30)$$

- **Multi-head Attention**

$$B^{\nu,\alpha',\alpha} = \frac{1}{\sqrt{d}} \sum_{\delta} K^{\nu,\delta,\alpha'} Q^{\nu,\delta,\alpha} \quad \left( B^{[n_h],[\ell],[\ell]} = \left[ \frac{1}{\sqrt{d}} (K^{\nu,[d],[\ell]})^T Q^{\nu,[d],[\ell]} \right]_{\nu \in [n_h]} \right) \quad (94.31)$$

$$A^{\nu,\delta,\alpha'} = \sum_{\alpha} V^{\nu,\delta,\alpha'} \underbrace{\text{softmax}(B^{\nu,[\ell],\alpha})(\alpha'|\alpha, \nu)}_{P(\alpha'|\nu,\alpha)} \quad (94.32)$$

$$\sum_{\alpha' \in [\ell]} P(\alpha' | \nu, \alpha) = 1 \quad (94.33)$$

$$A^{\nu, \delta, \alpha} \rightarrow A^{\Delta, \alpha} (A^{[n_h], [d], [\ell]} \rightarrow A^{[D], [\ell]}) \quad (94.34)$$

Column vector notation:

$$B^{\nu, \alpha', \alpha} = \frac{1}{\sqrt{d}} (K^{\nu, [d], \alpha'})^T Q^{\nu, [d], \alpha} \quad (94.35)$$

Important: Note that the softmax() makes the  $\alpha'$  component a probability, not the  $\alpha$  one!

For example, suppose  $\nu = 1$  (one head),  $\ell = 2$  (a 2 word segment), and  $d = 3$  (hidden dimension is 3). The  $Q^{[3], [2]}, K^{[3], [2]}, V^{[3], [2]}$  are  $3 \times 2$  matrices (i.e., two 3-dim column vectors). One uses the  $Q^{[3], [2]}$  and  $K^{[3], [2]}$  to arrive at a  $2 \times 2$  matrix  $P(\alpha' | \alpha)$  of probabilities. Then one uses that matrix of probabilities to replace

$$[V^{[3], 0}, V^{[3], 1}] \rightarrow [V^{[3], 0} P(0|0) + V^{[3], 1} P(1|0), V^{[3], 0} P(0|1) + V^{[3], 1} P(1|1)] \quad (94.36)$$

- **Positional Embedding Matrix**

$$E_{pos}^{[d], [\ell]}$$

$$E_{pos}^{\delta, \beta} = \begin{cases} \sin\left(\frac{\beta}{10^{4\delta/d}}\right) = \sin\left(2\pi \frac{\beta}{\lambda(\delta)}\right) & \text{if } \delta \text{ is even} \\ \cos\left(\frac{\beta}{10^{4(\delta-1)/d}}\right) = \cos\left(2\pi \frac{\beta}{\lambda(\delta)}\right) & \text{if } \delta \text{ is odd} \end{cases} \quad (94.37)$$

$E_{pos}^{\delta, \beta}$  changes in phase by  $\pi/2$  every time  $\delta$  changes by 1. Its wavelength  $\lambda$  is independent of  $\beta$ , but increases rapidly with  $\delta$ , from  $\lambda(\delta = 0) = 2\pi * 1$  to  $\lambda(\delta = d) = 2\pi * 10^4$ .

Total Embedding equals initial embedding plus positional embedding:  $E = E_0 + E_{pos}$

The purpose of positional embedding is to take  $e^{\beta, \alpha}$  to  $e^{\delta, \alpha} = \sum_{\beta} E_{pos}^{\delta, \beta} e^{\beta, \alpha}$  where  $e^{\delta, \alpha}$  changes quickly as  $\delta$  (i.e., position) changes.

- **ReLU**

For a tensor  $T$  of arbitrary shape

$$ReLU(T) = (T)_+ = \max(0, T) \quad (94.38)$$

max element-wise

- **Feed Forward neural net**

$$F(e^{[1],[\ell]}) = ReLU(e^{[1],[\ell]}W_1^{[\ell],[d]} + b_1^{[1],[d]})W_2^{[d],[\ell]} + b_1^{[1],[\ell]} \quad (94.39)$$

$$F(e^{[\ell]}) = W_2^{[\ell],[d]}ReLU(W_1^{[d],[\ell]}e^{[\ell]} + b_1^{[d]}) + b_1^{[\ell]} \quad (94.40)$$

- **Softmax**

softmax() takes a vector and returns a vector of probabilities of the same length

$$e^{[n]} \rightarrow P^{[n]} \quad (94.41)$$

where

$$P^\alpha = \frac{\exp(e^\alpha)}{\sum_{\alpha \in [n]} \exp(e^\alpha)} \quad \left( P^{[n]} = \frac{\exp(e^{[n]})}{\| \exp(e^{[n]}) \|_0} \right) \quad (94.42)$$

For example,

$$(1, 0, 0) \rightarrow (e, 1, 1)/norm \quad (94.43)$$

$$(10, 0, 0) \rightarrow (e^{10}, 1, 1)/norm \approx (1, 0, 0) \quad (94.44)$$

For any  $a \in \mathbb{R}$ ,

$$(a, a, a) \rightarrow (1, 1, 1)/3 \quad (94.45)$$

- **Skip Connection (Add & Normalize)**

A **skip connection** is when you split the input to a **filter** into two streams, one stream goes through the filter, the other doesn't. The one that doesn't is then merged with the output of the filter via a **add & normalize** node. The reason for making skip connections is that the signal exiting a filter is usually full of jumps and kinks. By merging that filter output with some of the filter input, one smooths out the filter output to some degree. This makes back-propagation differentiation better behaved.

The filter might be a Multi-Head Attention or a Feed Forward NN.

Add & Normalize just means  $(A + B)/norm$  where  $A$  and  $B$  are the two input signals and “norm” is some norm of  $A + B$  (for instance,  $\|A + B\|_2$ ).

Normalization keeps the signal from growing too big and saturating the signal entering components upstream. Normalization can also involve subtracting the mean  $\langle X \rangle$  of the signal  $X$  so as to get a signal  $X - \langle X \rangle$  that has zero mean.

- **Plates**

Fig.94.4 uses plates (a.k.a. stacks). These are explained in Chapter 71.

- **Redundancy**

For better results, the Decoder contains a pair of multi-head Attentions in series, and  $\Lambda$  of those pairs in parallel (i.e.,  $\Lambda$  layers of a plate).

- **Right Shifted Outputs**

“Outputs (Shifted Right)” in Fig.94.4 refers to what is called **forced teaching** in the RNN (recurrent neural net) literature. We explain forced teaching in Fig.94.5.

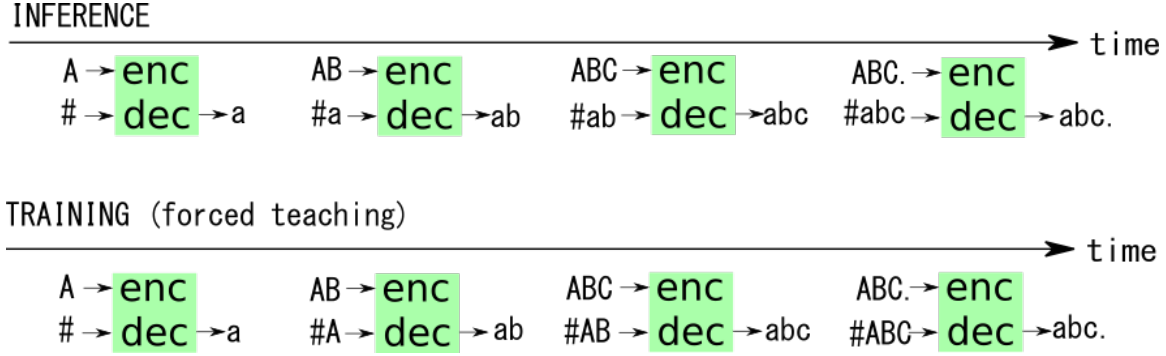


Figure 94.5: Training and Inference for vanilla transformer. “enc” and “dec” denote the encoder and decoder, respectively. A hash character represents the SOS (start of sentence) token, and a period represents the EOS (end of sentence) token. Capital letters represent ground truth tokens, and lower case ones represent predictions.

- **Masked Attention**

$$P(\alpha' | \nu, \alpha) = 0 \text{ if } \alpha' < \alpha \quad (94.46)$$

$\alpha$  and  $\alpha'$  are sentence positions and  $\alpha'$  is in the future (downstream) compared to  $\alpha$ . So as to not violate causality, this condition enforces the constraint that no attention is paid to sentence positions in the future  $\alpha$ .

- **plates**

Fig.94.4 uses plates (a.k.a. stacks). These are explained in Chapter 71.

### 94.2.1 Single Head Attention

Fig.94.6 gives a bnet representation of the “Single Head Attention” portion of Fig.94.4. The structural equations for that bnet, printed in blue, are as follows.

$$A^{[d],[\ell]} = V^{[d],[\ell]} P^{[\ell],[\ell]} \left( \text{Note that } \sum_{\alpha \in [\ell]} P^{\alpha, [\ell]} = 1 \right) \quad (94.47a)$$

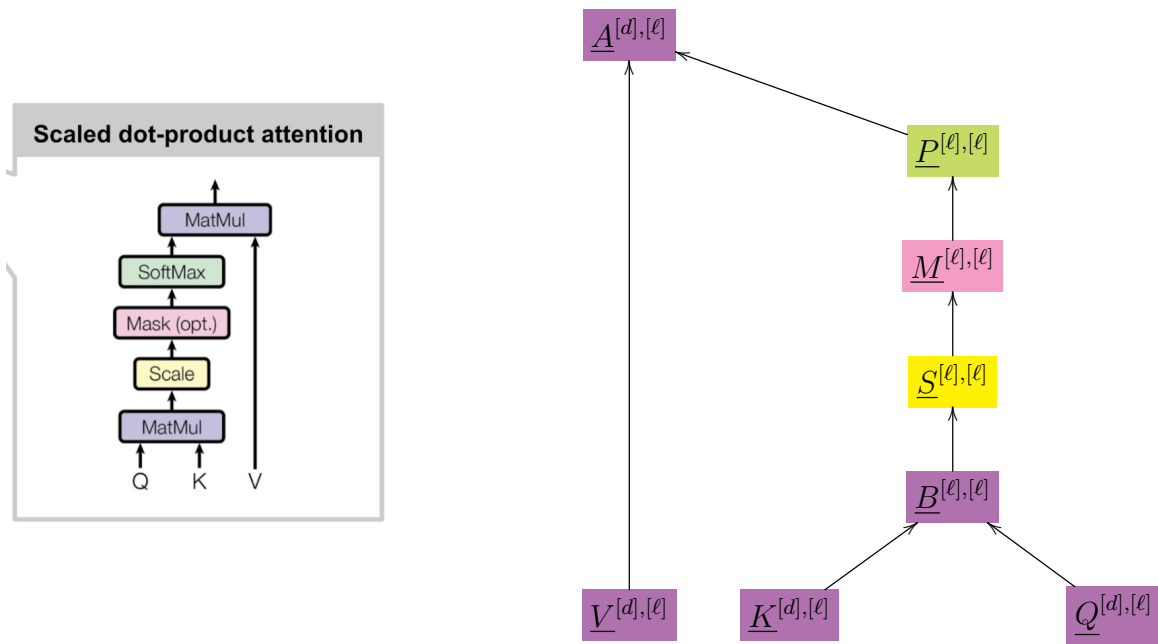


Figure 94.6: Single Head Attention. (Scaled Dot Product)

$$B^{[\ell],[\ell]} = (K^{[d],[\ell]})^T Q^{[d],[\ell]} \quad (94.47b)$$

$$K^{[d],[\ell]} = \text{prior} \quad (94.47c)$$

$$M^{[\ell],[\ell]} = \text{mask}(S^{[\ell],[\ell]}) \quad (94.47d)$$

$$P^{[\ell],[\ell]} = \text{softmax}(M^{[\ell],[\ell]}) \quad \left( \text{Note that } \sum_{\alpha \in [\ell]} P^{\alpha,[\ell]} = 1 \right) \quad (94.47e)$$

$$Q^{[d],[\ell]} = \text{prior} \quad (94.47f)$$

$$S^{[\ell],[\ell]} = \frac{B^{[\ell],[\ell]}}{\sqrt{d}} \quad (94.47g)$$



$$V^{[d],[\ell]} = \text{prior} \quad (94.47h)$$

## 94.2.2 Multi-Head Attention

Fig.94.7 gives a bnet representation of the “Multi-Head Attention” portion of Fig.94.4. The structural equations for that bnet, printed in blue, are as follows.

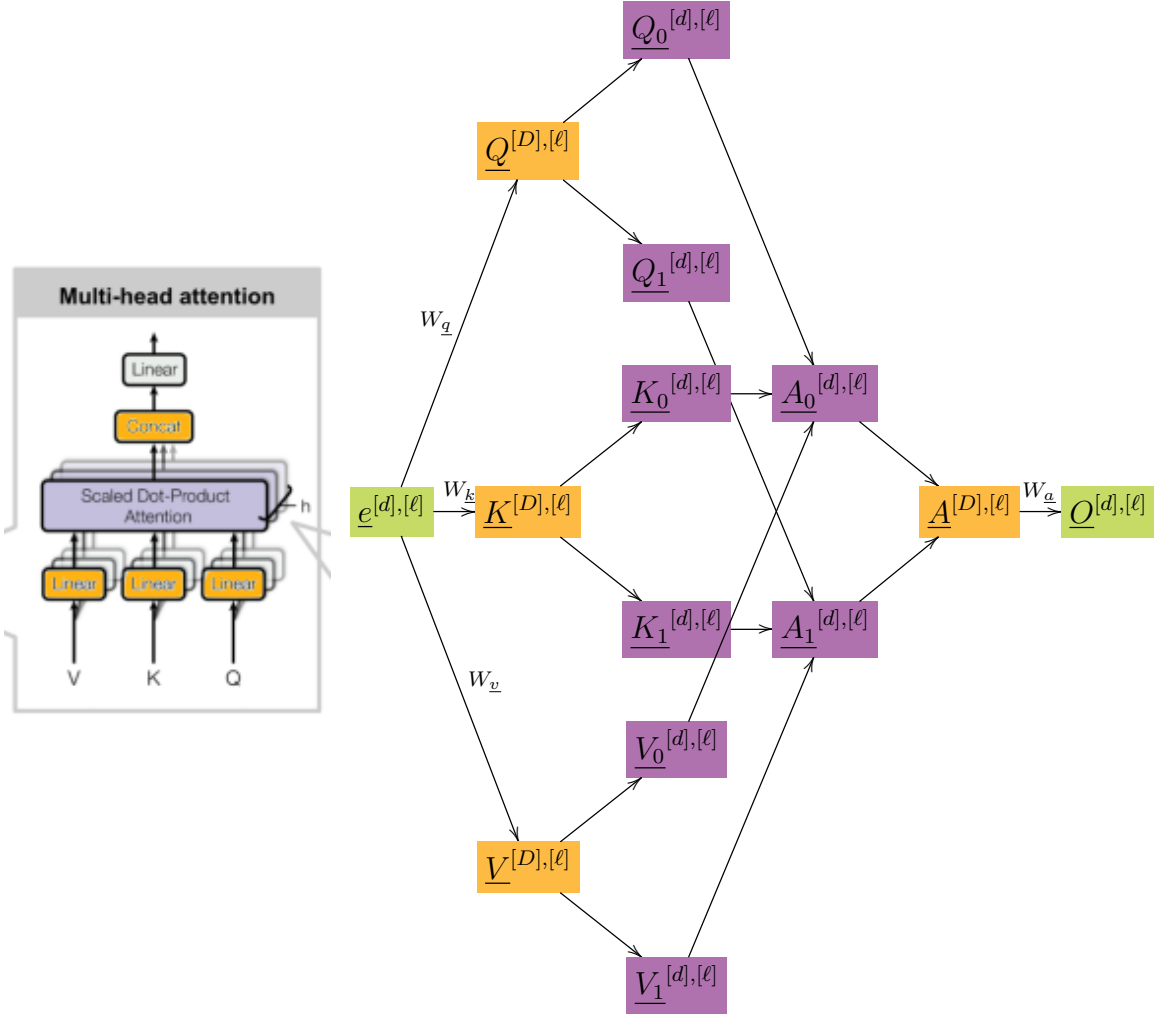


Figure 94.7: Multi-head Attention with 2 heads. Note that the orange nodes all have the same tensor shape.

$$A^{[D],[\ell]} = [A_0^{[d],[\ell]} | A_1^{[d],[\ell]}] \quad (94.48a)$$

$$A_0^{[d],[\ell]} = \text{Attention}(V_0^{[d],[\ell]}, K_0^{[d],[\ell]}, Q_0^{[d],[\ell]}) \quad (94.48b)$$

$$A_1^{[d],[\ell]} = \text{Attention}(V_1^{[d],[\ell]}, K_1^{[d],[\ell]}, Q_1^{[d],[\ell]}) \quad (94.48c)$$

$$K^{[D],[\ell]} = W_{\underline{k}}^{[D],[d]} e^{[d],[\ell]} \quad (94.48d)$$

$$K_0^{[d],[\ell]} = \text{linear}(K^{[D],[\ell]}) \text{ (split, then project a component)} \quad (94.48e)$$

$$K_1^{[d],[\ell]} = \text{linear}(K^{[D],[\ell]}) \text{ (split, then project a component)} \quad (94.48f)$$

$$O^{[d],[\ell]} = W_{\underline{a}}^{[d],[D]} A^{[D],[\ell]} \quad (94.48g)$$

$$Q^{[D],[\ell]} = W_{\underline{q}}^{[D],[d]} e^{[d],[\ell]} \quad (94.48h)$$

$$Q_0^{[d],[\ell]} = \text{linear}(Q^{[D],[\ell]}) \text{ (split, then project a component)} \quad (94.48i)$$

$$Q_1^{[d],[\ell]} = \text{linear}(Q^{[D],[\ell]}) \text{ (split, then project a component)} \quad (94.48j)$$

$$V^{[D],[\ell]} = W_{\underline{v}}^{[D],[d]} e^{[d],[\ell]} \quad (94.48k)$$

$$V_0^{[d],[\ell]} = \text{linear}(V^{[D],[\ell]}) \text{ (split, then project a component)} \quad (94.48l)$$

$$V_1^{[d],[\ell]} = \text{linear}(V^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (94.48\text{m})$$

$$e^{[d],[\ell]} = \text{prior} \quad (94.48\text{n})$$

### 94.2.3 Encoder

Fig.94.8 gives a bnet representation of the “Encoder” portion of Fig.94.4. The structural equations for that bnet, printed in blue, are as follows.

$$A^{[\Lambda],[D],[\ell]} = \text{Attention}(Q^{[\Lambda],[D],[\ell]}, K^{[\Lambda],[D],[\ell]}, V^{[\Lambda],[D],[\ell]}) \quad (94.49\text{a})$$

$$F^{[\Lambda],[d],[\ell]} = \text{feed\_forward\_nn}(N^{[\Lambda],[d],[\ell]}) \quad (94.49\text{b})$$

$$K^{[\Lambda],[D],[\ell]} = W_{\underline{k}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.49\text{c})$$

$$N^{[\Lambda],[d],[\ell]} = \text{normalize}(e^{[\Lambda],[d],[\ell]} + W_{\underline{a}}^{[d],[D]} A^{[\Lambda],[D],[\ell]}) \quad (94.49\text{d})$$

$$Q^{[\Lambda],[D],[\ell]} = W_{\underline{q}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.49\text{e})$$

$$V^{[\Lambda],[D],[\ell]} = W_{\underline{v}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.49\text{f})$$

$$e^{[\Lambda],[d],[\ell]} = E^{[\Lambda],[d],[L]} x^{[L],[\ell]} \quad (94.49\text{g})$$

$$n^{[\Lambda],[d],[\ell]} = \text{normalize}(N^{[\Lambda],[d],[\ell]} + F^{[\Lambda],[d],[\ell]}) \quad (94.49\text{h})$$

$$x^{[L],[\ell]} = \text{prior} \quad (94.49\text{i})$$

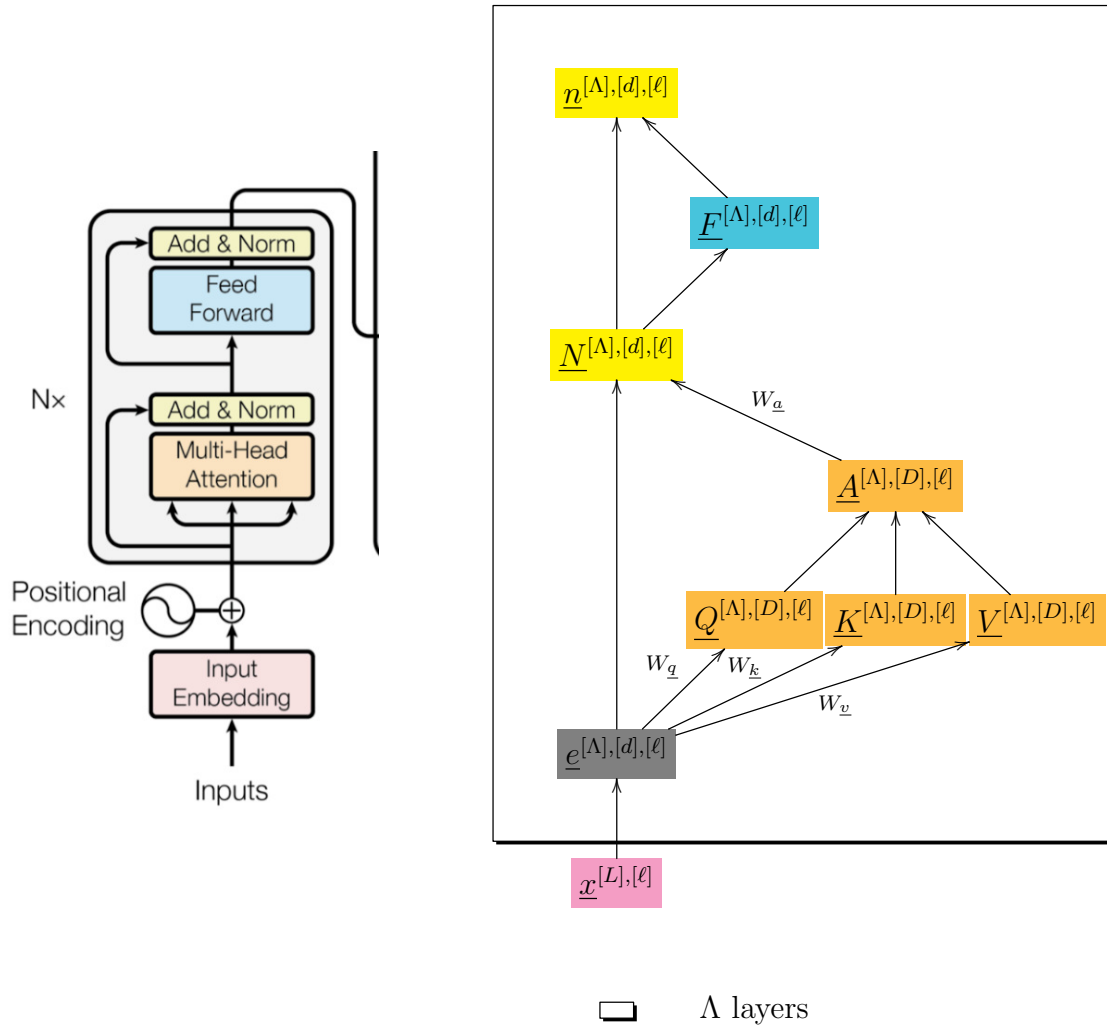


Figure 94.8: Encoder of Vanilla Transformer Net.

#### 94.2.4 Decoder

Fig.94.9 gives a bnet representation of the “Decoder” portion of Fig.94.4. The structural equations for that bnet, printed in blue, are as follows.

$$A^{[\Lambda],[D],[\ell]} = \text{Attention}(Q^{[\Lambda],[D],[\ell]}, K^{[\Lambda],[D],[\ell]}, V^{[\Lambda],[D],[\ell]}) \quad (94.50a)$$

$$F^{[\Lambda],[d],[\ell]} = \text{feed\_forward\_nn}(j^{[\Lambda],[d],[\ell]}) \quad (94.50b)$$

$$I^{[L],[\ell]} = W^{[L],[\Lambda],[d]} Y^{[\Lambda],[d],[\ell]} \quad (94.50c)$$

$$J^{[\Lambda],[d],[\ell]} = \text{normalize}(W_{\underline{a}}^{[d],[D]} A^{[\Lambda],[D],[\ell]} + e^{[\Lambda],[d],[\ell]}) \quad (94.50d)$$

$$K^{[\Lambda],[D],[\ell]} = W_{\underline{k}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.50e)$$

$$P^{[L],[\ell]} = \text{softmax}(I^{[L],[\ell]}) \quad (\sum_{\alpha \in [\ell]} P^{[L],\alpha} = 1) \quad (94.50f)$$

$$Q^{[\Lambda],[D],[\ell]} = W_{\underline{q}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.50g)$$

$$V^{[\Lambda],[D],[\ell]} = W_{\underline{v}}^{[D],[d]} e^{[\Lambda],[d],[\ell]} \quad (94.50h)$$

$$Y^{[\Lambda],[d],[\ell]} = \text{normalize}(F^{[\Lambda],[d],[\ell]} + J^{[\Lambda],[d],[\ell]}) \quad (94.50i)$$

$$a^{[\Lambda],[D],[\ell]} = \text{Attention}(v^{[\Lambda],[D],[\ell]}, k^{[\Lambda],[D],[\ell]}, q^{[\Lambda],[D],[\ell]}) \quad (94.50j)$$

$$e^{[\Lambda],[d],[\ell]} = E^{[\Lambda],[d],[L]} x^{[L],[\ell]} \quad (94.50k)$$

$$j^{[\Lambda],[d],[\ell]} = \text{normalize}(U_{\underline{a}}^{[d],[D]} a^{[\Lambda],[D],[\ell]} + J^{[\Lambda],[d],[\ell]}) \quad (94.50l)$$

$$k^{[\Lambda],[D],[\ell]} = U_{\underline{k}}^{[D],[d]} n^{[\Lambda],[d],[\ell]} \quad (94.50m)$$

$$n^{[\Lambda],[d],[\ell]} = \text{Prior coming from Encoder.} \quad (94.50n)$$

$$q^{[\Lambda],[D],[\ell]} = U_{\underline{q}}^{[D],[d]} J^{[\Lambda],[d],[\ell]} \quad (94.50o)$$

$$v^{[\Lambda],[D],[\ell]} = U_{\underline{v}}^{[D],[d]} n^{[\Lambda],[d],[\ell]} \tag{94.50p}$$

$$x^{[L],[\ell]} = \text{prior, right shifted output} \tag{94.50q}$$

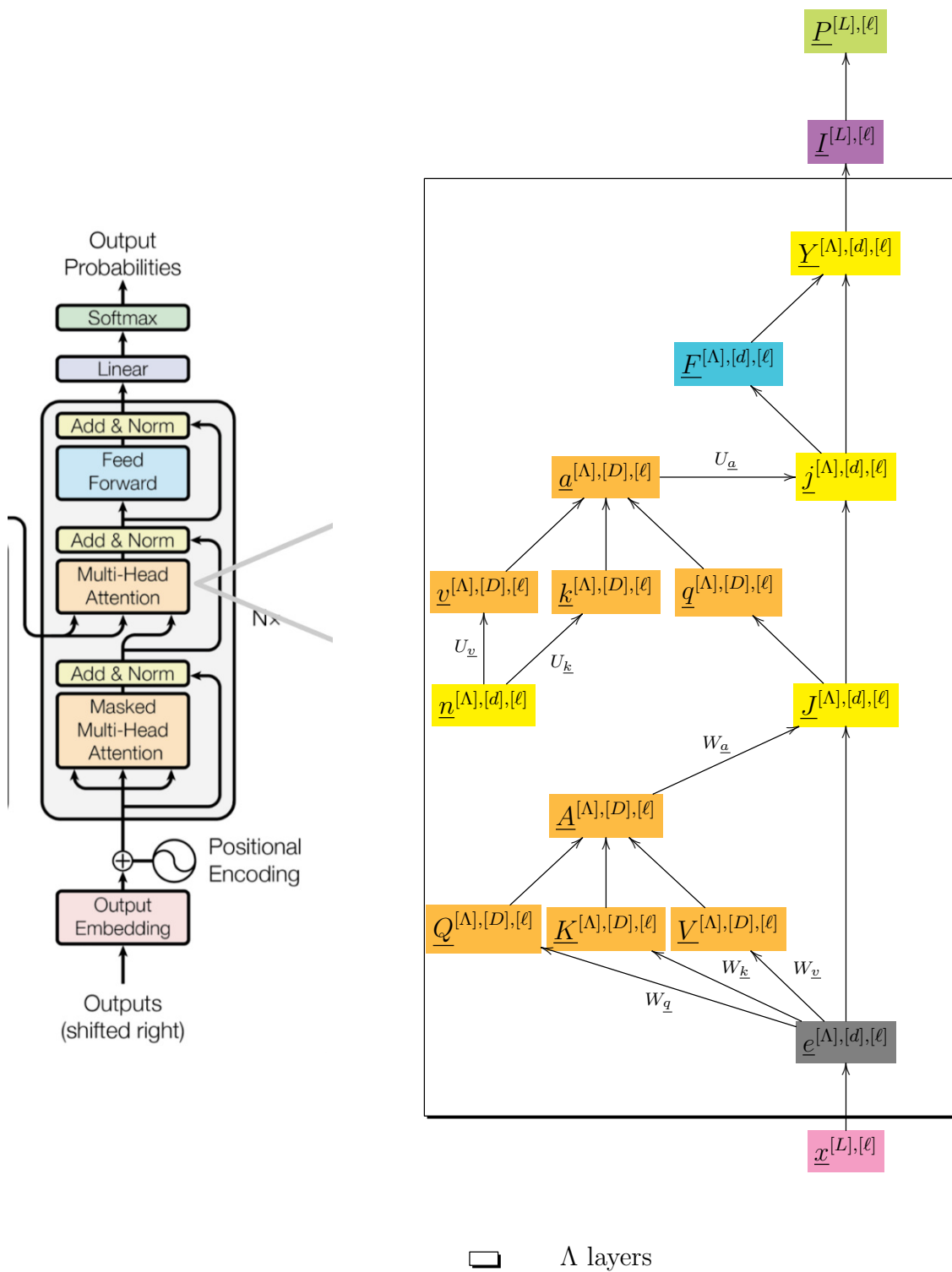


Figure 94.9: Decoder of Vanilla Transformer Net.

# Bibliography

- [1] Alan Agresti. *An introduction to categorical data analysis*. John Wiley & Sons, 2018.
- [2] Elias Bareinboim, Jin Tian, and Judea Pearl. Recovering from selection bias in causal and statistical inference. In *Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r425.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r425.pdf).
- [3] Dan Bendel. Metropolis-Hastings: A comprehensive overview and proof. <https://similarweb.engineering/mcmc/>.
- [4] David Benkeser and Antoine Chambaz. A ride in targeted learning territory. <https://achambaz.github.io/tlride/tlride-book.pdf>.
- [5] Alexandra M Carvalho. Scoring functions for learning Bayesian networks. [http://www.lx.it.pt/~asmc/pub/talks/09-TA/ta\\_pres.pdf](http://www.lx.it.pt/~asmc/pub/talks/09-TA/ta_pres.pdf).
- [6] Bo Chang. Copula: a very short introduction, article in Bo’s Blog. <https://bochang.me/blog/posts/copula/>.
- [7] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. <https://arxiv.org/abs/1603.02754>.
- [8] Victor Chernozhukov, Carlos Cinelli, Whitney Newey, Amit Sharma, and Vasilis Syrgkanis. Long story short: Omitted variable bias in causal machine learning. [https://www.nber.org/system/files/working\\_papers/w30302/w30302.pdf](https://www.nber.org/system/files/working_papers/w30302/w30302.pdf).
- [9] Carlos Cinelli, Andrew Forney, and Judea Pearl. A crash course in good and bad controls. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r493.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r493.pdf).
- [10] Carlos Cinelli and Chad Hazlett. Making sense of sensitivity: Extending omitted variable bias. [https://carloscinelli.com/files/Cinelli%20and%20Hazlett%20\(2020\)%20-%20Making%20Sense%20of%20Sensitivity.pdf](https://carloscinelli.com/files/Cinelli%20and%20Hazlett%20(2020)%20-%20Making%20Sense%20of%20Sensitivity.pdf).
- [11] Scott Cunningham. *Causal inference: The mixtape*. Yale University Press, 2021. <https://mixtape.scunning.com/index.html>.



- [12] Robin J. Evans. Graphical methods for inequality constraints in marginalized DAGs. <https://arxiv.org/abs/1209.2978>.
- [13] Matheus Facure Alves. *Causal Inference for The Brave and True*. 2021. <https://matheusfacure.github.io/python-causality-handbook/landing-page.html>.
- [14] George Fei. Modeling uplift directly: Uplift decision tree with kl divergence and euclidean distance as splitting criteria. <https://www.aboutwayfair.com/tech-innovation/modeling-uplift-directly-uplift-decision-tree-with-kl-divergence-and-euclidean-distance-as-splitting-criteria>.
- [15] Charles Fox, Neil Girdhar, and Kevin Gurney. A causal Bayesian network view of reinforcement learning. <https://www.aaai.org/Papers/FLAIRS/2008/FLAIRS08-030.pdf>.
- [16] Bruno Gonçalves. Model testing and causal search. blog post <https://medium.com/data-for-science/causal-inference-part-vii-model-testing-and-causal-search-536b796f0384>.
- [17] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, David Warde-Farley Bing Xu, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. <https://arxiv.org/abs/1406.2661>.
- [18] Pierre Gutierrez and Jean-Yves Gérardy. Causal inference and uplift modelling: A review of the literature. In *Proceedings of The 3rd International Conference on Predictive Applications and APIs*, pages 1–13, 2017. <http://proceedings.mlr.press/v67/gutierrez17a.html>.
- [19] James Douglas Hamilton. *Time series analysis*. Princeton University Press, 2020.
- [20] Sebastian Haneuse and Andrea Rotnitzky. Estimation of the effect of interventions that modify the received treatment. *Statistics in medicine*, 32(30):5260–5277, 2013. Main:<https://sci-hub.se/10.1002/sim.5907>, Supplement: <https://onlinelibrary.wiley.com/action/downloadSupplement?doi=10.1002%2Fsim.5907&file=sim5907-sup-0001-Appendix.pdf>.
- [21] Christina Heinze-Deml. Causality, spring semester 2019 at ETH Zurich. [https://stat.ethz.ch/lectures/ss19/causality.php#course\\_materials](https://stat.ethz.ch/lectures/ss19/causality.php#course_materials).
- [22] MA Hernán and J Robins. Causal inference: What if. Boca Raton: Chapman & Hill/CRC, <https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>.

- [23] Katherine Hoffman. An illustrated guide to modified treatment policies, part 1: Introduction and motivation, article in KHstats blog. <https://www.khstats.com/blog/lmtp/lmtp/>.
- [24] Katherine Hoffman. An illustrated guide to TMLE, article in KHstats blog. <https://www.khstats.com/blog/tmle/tutorial/>.
- [25] Cecil Huang and Adnan Darwiche. Inference in belief networks: A procedural guide. *International journal of approximate reasoning*, 15(3):225–263, 1996. <http://www.ar-tiste.com/Huang-Darwiche1996.pdf>.
- [26] Tommi S. Jaakkola and Michael I. Jordan. Variational probabilistic inference and the QMR-DT network. <http://arxiv.org/abs/1105.5462>.
- [27] Michael I. Jordan. course: Stat260-Bayesian modeling and inference, lecture date: February 8-2010, title: The conjugate prior for the Normal distribution. <https://people.eecs.berkeley.edu/~jordan/courses/260-spring10/lectures/lecture5.pdf>.
- [28] Chaitanya K. Joshi. Transformer (machine learning model). <https://graphdeeplearning.github.io/post/transformers-are-gnns/>.
- [29] Chung-Ming Kuan. Introduction to time series analysis, Fall 2014 lectures given at the Department of Finance, National Taiwan University. [https://homepage.ntu.edu.tw/~ckuan/pdf/2014fall/Lec-TimeSeries\\_slide-Fall2014.pdf](https://homepage.ntu.edu.tw/~ckuan/pdf/2014fall/Lec-TimeSeries_slide-Fall2014.pdf).
- [30] Steffen L Lauritzen and David J Spiegelhalter. Local computations with probabilities on graphical structures and their application to expert systems. *Journal of the Royal Statistical Society: Series B (Methodological)*, 50(2):157–194, 1988. <http://www.eecis.udel.edu/~shatkay/Course/papers/Lauritzen1988.pdf>.
- [31] Sergey Levine. Course CS 285 at UC Berkeley, Deep reinforcement learning. <http://rail.eecs.berkeley.edu/deeprlcourse/>.
- [32] Ang Li. Ph.D. Thesis, UCLA 2021. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r507.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r507.pdf).
- [33] Dimitris Margaritis. Learning Bayesian network model structure from data (thesis, 2003, Carnegie Mellon Univ). <https://apps.dtic.mil/sti/citations/ADA461103>.
- [34] Adam A Margolin, Ilya Nemenman, Katia Basso, Chris Wiggins, Gustavo Stolovitzky, Riccardo Dalla Favera, and Andrea Califano. Aracne: an algorithm for the reconstruction of gene regulatory networks in a mammalian

- cellular context. In *BMC bioinformatics*, volume 7, page S7. Springer, 2006. <https://link.springer.com/article/10.1186/1471-2105-7-S1-S7>.
- [35] Samuele Mazzanti. Black-box models are actually more explainable than a logistic regression. <https://towardsdatascience.com/black-box-models-are-actually-more-explainable-than-a-logistic-regression-f263c22795d>.
  - [36] Samuele Mazzanti. SHAP values explained exactly how you wished someone explained to you. <https://towardsdatascience.com/shap-explained-the-way-i-wish-someone-explained-it-to-me-ab81cc69ef30>.
  - [37] Robert J. McEliece, David J. C. MacKay, and Jung-Fu Cheng. Turbo decoding as an instance of Pearl’s belief propagation algorithm. <http://authors.library.caltech.edu/6938/1/MCEieeejstc98.pdf>.
  - [38] Scott Mueller, Ang Li, and Judea Pearl. Causes of effects: Learning individual responses from population data. *arXiv preprint arXiv:2104.13730*, 2021. <https://arxiv.org/abs/2104.13730>.
  - [39] Brady Neal. Introduction to causal inference, Fall 2020, lectures and book. <https://www.bradyneal.com/causal-inference-course>.
  - [40] Richard E Neapolitan. *Learning Bayesian networks*. Pearson Prentice Hall, 2004.
  - [41] Andrew Ng. Lecture at deeplearning.ai on recurrent neural networks. <http://www.ar-tiste.com/ng-lec-rnn.pdf>.
  - [42] Gregory Nuel. Tutorial on exact belief propagation in Bayesian networks: from messages to algorithms. <https://arxiv.org/abs/1201.4724>.
  - [43] paperspace.com. PyTorch 101, Part 1: Understanding graphs, automatic differentiation and autograd. <https://blog.paperspace.com/pytorch-101-understanding-graphs-and-automatic-differentiation/>.
  - [44] Judea Pearl. Linear models: A useful microscope for causal analysis. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r409-corrected-reprint.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r409-corrected-reprint.pdf).
  - [45] Judea Pearl. Mediating instrumental variables. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r210.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r210.pdf).
  - [46] Judea Pearl. On the testability of causal models with latent and instrumental variables. <https://arxiv.org/abs/1302.4976>.

- [47] Judea Pearl. Reverend Bayes on inference engines: A distributed hierarchical approach. <https://www.aaai.org/Papers/AAAI/1982/AAAI82-032.pdf>, 1982.
- [48] Judea Pearl. *Probabilistic Inference in Intelligent Systems*. Morgan Kaufmann, 1988.
- [49] Judea Pearl. The causal mediation formula—a guide to the assessment of pathways and mechanisms. *Prevention science*, 13(4):426–436, 2012. <https://apps.dtic.mil/sti/pdfs/ADA557663.pdf>.
- [50] Judea Pearl. *Causality: Models, Reasoning, and Inference, Second Edition*. Cambridge University Press, 2013.
- [51] Judea Pearl. Causal and counterfactual inference. *The Handbook of Rationality*, pages 1–41, 2019. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r485.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r485.pdf).
- [52] Judea Pearl and Elias Bareinboim. Transportability of causal and statistical relations: A formal approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, 2011. <https://ojs.aaai.org/index.php/AAAI/article/view/7861>.
- [53] Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.
- [54] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic Books, 2018.
- [55] Judea Pearl and James M Robins. Probabilistic evaluation of sequential plans from causal models with hidden variables. *arXiv preprint arXiv:1302.4977*, 2013. <https://arxiv.org/abs/1302.4977>.
- [56] Ashwin Rao and Tikhon Jelvis. *Foundations of Reinforcement Learning with Applications in Finance*. <https://stanford.edu/~ashlearn/RLForFinanceBook/book.pdf>.
- [57] ReliaSoft. System analysis reference. [http://reliawiki.org/index.php/System\\_Analysis\\_Reference](http://reliawiki.org/index.php/System_Analysis_Reference).
- [58] Piotr Rzepakowski and Szymon Jaroszewicz. Decision trees for uplift modeling with single and multiple treatments. *Knowledge and Information Systems*, 32(2):303–327, 2012. <https://link.springer.com/content/pdf/10.1007/s10115-011-0434-0.pdf>.
- [59] Scholarpedia. Granger causality. <http://www.scholarpedia.org/article/Granger-causality>.

- [60] Marco Scutari. bnlearn. <https://www.bnlearn.com/>.
- [61] Marco Scutari, Catharina Elisabeth Graafland, and José Manuel Gutiérrez. Who learns better Bayesian network structures: Accuracy and speed of structure learning algorithms. *International Journal of Approximate Reasoning*, 115:235–253, 2019. <https://arxiv.org/abs/1805.11908>.
- [62] Nitish Srivastava, G E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. <http://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf>.
- [63] StatQuest. XGBoost Parts 1-4 (videos). <https://statquest.org/video-index/>.
- [64] Russell Stuart and Norvig Peter. Artificial intelligence-a modern approach 3rd ed, 2016.
- [65] Masayoshi Takahashi. Statistical inference in missing data by MCMC and non-MCMC multiple imputation algorithms: Assessing the effects of between-imputation iterations. *Data Science Journal*, 16, 2017. <https://datascience.codata.org/articles/10.5334/dsj-2017-037/>.
- [66] theinvestorsbook.com. Pert analysis. <https://theinvestorsbook.com/pert-analysis.html>.
- [67] Jin Tian and Judea Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics and Artificial Intelligence*, 28(1):287–313, 2000. [https://ftp.cs.ucla.edu/pub/stat\\_ser/r271-A.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r271-A.pdf).
- [68] Robert R. Tucci. Bayesian networks (a.k.a. causal models, DAGs) and the passage of time. blog post in blog Quantum Bayesian Networks, <https://qbnets.wordpress.com/2021/07/16/bayesian-networks-aka-causal-models-dags-and-the-passage-of-time/>.
- [69] Robert R. Tucci. Bell’s inequalities for Bayesian statisticians. blog post in blog Quantum Bayesian Networks, <https://qbnets.wordpress.com/2008/09/19/bells-inequities-for-bayesian-statistician/>.
- [70] Robert R. Tucci. Goodness of causal fit. [https://github.com/rrtucci/DAG\\_Lie\\_Detector](https://github.com/rrtucci/DAG_Lie_Detector).
- [71] Robert R. Tucci. JudeasRx. <https://github.com/rrtucci/JudeasRx>.
- [72] Robert R. Tucci. Mappa Mundi. [https://github.com/rrtucci/mappa\\_mundi](https://github.com/rrtucci/mappa_mundi).
- [73] Robert R. Tucci. Quantum d-separation and quantum belief propagation. <https://arxiv.org/abs/2012.09635>.

- [74] Robert R. Tucci. Quantum Fog. <https://github.com/artiste-qb-net/quantum-fog>.
- [75] Robert R. Tucci. SCuMpy. <https://github.com/rrtucci/scumpy>.
- [76] Robert R. Tucci. Shannon information theory without shedding tears over delta & epsilon proofs or typical sequences. <https://arxiv.org/abs/1208.2737>.
- [77] Robert R. Tucci. texnn. <https://github.com/rrtucci/texnn>.
- [78] Ashish Vaswani et al. Attention is all you need. <https://arxiv.org/abs/1706.03762>.
- [79] W. E. Vesely, F. F. Goldberg, N. H. Roberts, and D. F. Haasl. Fault tree handbook nureg-0492. <https://www.nrc.gov/reading-rm/doc-collections/nuregs/staff/sr0492/>.
- [80] Martin J Wainwright and Michael I Jordan. Graphical models, exponential families, and variational inference (book). <https://people.eecs.berkeley.edu/~wainwrig/Papers/WaiJor08.FTML.pdf>.
- [81] Lilian Weng. The multi-armed bandit problem and its solutions. [lilianweng.github.io/lil-log](http://lilianweng.github.io/lil-log), <http://lilianweng.github.io/lil-log/2018/01/23/the-multi-armed-bandit-problem-and-its-solutions.html>, 2018.
- [82] Lilian Weng. What are diffusion models? [lilianweng.github.io/lil-log](http://lilianweng.github.io/lil-log), <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>, 2021.
- [83] Wikibooks. Control systems. [https://en.wikibooks.org/wiki/Control\\_Systems](https://en.wikibooks.org/wiki/Control_Systems).
- [84] Wikipedia. AdaBoost. <https://en.wikipedia.org/wiki/AdaBoost>.
- [85] Wikipedia. Autoregressive model. [https://en.wikipedia.org/wiki/Autoregressive\\_model](https://en.wikipedia.org/wiki/Autoregressive_model).
- [86] Wikipedia. Autoregressive moving-average model. [https://en.wikipedia.org/wiki/Autoregressive%E2%80%93moving-average\\_model](https://en.wikipedia.org/wiki/Autoregressive%E2%80%93moving-average_model).
- [87] Wikipedia. Baum-Welsh algorithm. [https://en.wikipedia.org/wiki/Baum%E2%80%93Welch\\_algorithm](https://en.wikipedia.org/wiki/Baum%E2%80%93Welch_algorithm).
- [88] Wikipedia. Belief propagation. [https://en.wikipedia.org/wiki/Belief\\_propagation](https://en.wikipedia.org/wiki/Belief_propagation).
- [89] Wikipedia. Berkson’s paradox. [https://en.wikipedia.org/wiki/Berkson%27s\\_paradox](https://en.wikipedia.org/wiki/Berkson%27s_paradox).

- [90] Wikipedia. Bernoulli distribution. [https://en.wikipedia.org/wiki/Bernoulli\\_distribution](https://en.wikipedia.org/wiki/Bernoulli_distribution).
- [91] Wikipedia. Beta distribution. [https://en.wikipedia.org/wiki/Beta\\_distribution](https://en.wikipedia.org/wiki/Beta_distribution).
- [92] Wikipedia. Beta function. [https://en.wikipedia.org/wiki/Beta\\_function](https://en.wikipedia.org/wiki/Beta_function).
- [93] Wikipedia. Binary decision diagram. [https://en.wikipedia.org/wiki/Binary\\_decision\\_diagram](https://en.wikipedia.org/wiki/Binary_decision_diagram).
- [94] Wikipedia. Boolean algebra. [https://en.wikipedia.org/wiki/Boolean\\_algebra](https://en.wikipedia.org/wiki/Boolean_algebra).
- [95] Wikipedia. Bootstrap aggregating. [https://en.wikipedia.org/wiki/Bootstrap\\_aggregating](https://en.wikipedia.org/wiki/Bootstrap_aggregating).
- [96] Wikipedia. Categorical distribution. [https://en.wikipedia.org/wiki/Categorical\\_distribution](https://en.wikipedia.org/wiki/Categorical_distribution).
- [97] Wikipedia. Chi-square distribution. [https://en.wikipedia.org/wiki/Chi-square\\_distribution](https://en.wikipedia.org/wiki/Chi-square_distribution).
- [98] Wikipedia. Chow-Liu tree. [https://en.wikipedia.org/wiki/Chow%E2%80%933Liu\\_tree](https://en.wikipedia.org/wiki/Chow%E2%80%933Liu_tree).
- [99] Wikipedia. Cochran's theorem. [https://en.wikipedia.org/wiki/Cochran%E2%80%9327s\\_theorem](https://en.wikipedia.org/wiki/Cochran%E2%80%9327s_theorem).
- [100] Wikipedia. Conjugate prior. [https://en.wikipedia.org/wiki/Conjugate\\_prior](https://en.wikipedia.org/wiki/Conjugate_prior).
- [101] Wikipedia. Copula. [https://en.wikipedia.org/wiki/Copula\\_\(probability\\_theory\)](https://en.wikipedia.org/wiki/Copula_(probability_theory)).
- [102] Wikipedia. Cramer-Rao bound. [https://en.wikipedia.org/wiki/Cram%C3%A9r%E2%80%93Rao\\_bound](https://en.wikipedia.org/wiki/Cram%C3%A9r%E2%80%93Rao_bound).
- [103] Wikipedia. Cross-validation. [https://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics)).
- [104] Wikipedia. Data processing inequality. [https://en.wikipedia.org/wiki/Data\\_processing\\_inequality](https://en.wikipedia.org/wiki/Data_processing_inequality).
- [105] Wikipedia. Dirichlet distribution. [https://en.wikipedia.org/wiki/Dirichlet\\_distribution](https://en.wikipedia.org/wiki/Dirichlet_distribution).

- [106] Wikipedia. Errors in variables models. [https://en.wikipedia.org/wiki/Errors-in-variables\\_models](https://en.wikipedia.org/wiki/Errors-in-variables_models).
- [107] Wikipedia. Expectation maximization. [https://en.wikipedia.org/wiki/Expectation%E2%80%93maximization\\_algorithm](https://en.wikipedia.org/wiki/Expectation%E2%80%93maximization_algorithm).
- [108] Wikipedia. F-distribution. <https://en.wikipedia.org/wiki/F-distribution>.
- [109] Wikipedia. Frisch-Waugh-Lovell theorem. [https://en.wikipedia.org/wiki/Frisch%E2%80%93Waugh%E2%80%93Lovell\\_theorem](https://en.wikipedia.org/wiki/Frisch%E2%80%93Waugh%E2%80%93Lovell_theorem).
- [110] Wikipedia. Functional derivative. [https://en.wikipedia.org/wiki/Functional\\_derivative](https://en.wikipedia.org/wiki/Functional_derivative).
- [111] Wikipedia. Gamma distribution. [https://en.wikipedia.org/wiki/Gamma\\_distribution](https://en.wikipedia.org/wiki/Gamma_distribution).
- [112] Wikipedia. Gamma function. [https://en.wikipedia.org/wiki/Gamma\\_function](https://en.wikipedia.org/wiki/Gamma_function).
- [113] Wikipedia. Gated recurrent unit. [https://en.wikipedia.org/wiki/Gated\\_recurrent\\_unit](https://en.wikipedia.org/wiki/Gated_recurrent_unit).
- [114] Wikipedia. Gibbs sampling. [https://en.wikipedia.org/wiki/Gibbs\\_sampling](https://en.wikipedia.org/wiki/Gibbs_sampling).
- [115] Wikipedia. Granger causality. [https://en.wikipedia.org/wiki/Granger\\_causality](https://en.wikipedia.org/wiki/Granger_causality).
- [116] Wikipedia. Hidden Markov model. [https://en.wikipedia.org/wiki/Hidden\\_Markov\\_model](https://en.wikipedia.org/wiki/Hidden_Markov_model).
- [117] Wikipedia. Hoeffding's inequality. [https://en.wikipedia.org/wiki/Hoeffding%27s\\_inequality](https://en.wikipedia.org/wiki/Hoeffding%27s_inequality).
- [118] Wikipedia. Importance sampling. [https://en.wikipedia.org/wiki/Importance\\_sampling](https://en.wikipedia.org/wiki/Importance_sampling).
- [119] Wikipedia. Instrumental variables estimation. [https://en.wikipedia.org/wiki/Instrumental\\_variables\\_estimation](https://en.wikipedia.org/wiki/Instrumental_variables_estimation).
- [120] Wikipedia. Inverse transform sampling. [https://en.wikipedia.org/wiki/Inverse\\_transform\\_sampling](https://en.wikipedia.org/wiki/Inverse_transform_sampling).
- [121] Wikipedia. Jackknife resampling. [https://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics)).



- [122] Wikipedia. Junction tree algorithm. [https://en.wikipedia.org/wiki/Junction\\_tree\\_algorithm](https://en.wikipedia.org/wiki/Junction_tree_algorithm).
- [123] Wikipedia. k-means clustering. [https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering).
- [124] Wikipedia. Kalman filter. [https://en.wikipedia.org/wiki/Kalman\\_filter](https://en.wikipedia.org/wiki/Kalman_filter).
- [125] Wikipedia. Kernel method. [https://en.wikipedia.org/wiki/Kernel\\_method](https://en.wikipedia.org/wiki/Kernel_method).
- [126] Wikipedia. Kernel perceptron. [https://en.wikipedia.org/wiki/Kernel\\_perceptron](https://en.wikipedia.org/wiki/Kernel_perceptron).
- [127] Wikipedia. Laplace transform. [https://en.wikipedia.org/wiki/Laplace\\_transform](https://en.wikipedia.org/wiki/Laplace_transform).
- [128] Wikipedia. Least squares. [https://en.wikipedia.org/wiki/Least\\_squares](https://en.wikipedia.org/wiki/Least_squares).
- [129] Wikipedia. Legendre transformation. [https://en.wikipedia.org/wiki/Legendre\\_transformation](https://en.wikipedia.org/wiki/Legendre_transformation).
- [130] Wikipedia. Likelihood-ratio test. [https://en.wikipedia.org/wiki/Likelihood-ratio\\_test](https://en.wikipedia.org/wiki/Likelihood-ratio_test).
- [131] Wikipedia. Linear regression. [https://en.wikipedia.org/wiki/Linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression).
- [132] Wikipedia. Long short term memory. [https://en.wikipedia.org/wiki/Long\\_short-term\\_memory](https://en.wikipedia.org/wiki/Long_short-term_memory).
- [133] Wikipedia. Markov blanket. [https://en.wikipedia.org/wiki/Markov\\_blanket](https://en.wikipedia.org/wiki/Markov_blanket).
- [134] Wikipedia. Metropolis-Hastings method. [https://en.wikipedia.org/wiki/Metropolis%E2%80%93Hastings\\_algorithm](https://en.wikipedia.org/wiki/Metropolis%E2%80%93Hastings_algorithm).
- [135] Wikipedia. Minimum spanning tree. [https://en.wikipedia.org/wiki/Minimum\\_spanning\\_tree](https://en.wikipedia.org/wiki/Minimum_spanning_tree).
- [136] Wikipedia. Monte Carlo methods. [https://en.wikipedia.org/wiki/Category:Monte\\_Carlo\\_methods](https://en.wikipedia.org/wiki/Category:Monte_Carlo_methods).
- [137] Wikipedia. Moving-average model. [https://en.wikipedia.org/wiki/Moving-average\\_model](https://en.wikipedia.org/wiki/Moving-average_model).
- [138] Wikipedia. Multinomial distribution. [https://en.wikipedia.org/wiki/Multinomial\\_distribution](https://en.wikipedia.org/wiki/Multinomial_distribution).

- [139] Wikipedia. Multinomial theorem. [https://en.wikipedia.org/wiki/Multinomial\\_theorem](https://en.wikipedia.org/wiki/Multinomial_theorem).
- [140] Wikipedia. Multivariate normal distribution. [https://en.wikipedia.org/wiki/Multivariate\\_normal\\_distribution](https://en.wikipedia.org/wiki/Multivariate_normal_distribution).
- [141] Wikipedia. Natural experiment. [https://en.wikipedia.org/wiki/Natural\\_experiment](https://en.wikipedia.org/wiki/Natural_experiment).
- [142] Wikipedia. Non-negative matrix factorization. [https://en.wikipedia.org/wiki/Non-negative\\_matrix\\_factorization](https://en.wikipedia.org/wiki/Non-negative_matrix_factorization).
- [143] Wikipedia. Ordinary least squares. [https://en.wikipedia.org/wiki/Ordinary\\_least\\_squares](https://en.wikipedia.org/wiki/Ordinary_least_squares).
- [144] Wikipedia. Program evaluation and review technique. [https://en.wikipedia.org/wiki/Program\\_evaluation\\_and\\_review\\_technique](https://en.wikipedia.org/wiki/Program_evaluation_and_review_technique).
- [145] Wikipedia. Random forest. [https://en.wikipedia.org/wiki/Random\\_forest](https://en.wikipedia.org/wiki/Random_forest).
- [146] Wikipedia. Receiver operating characteristic. [https://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](https://en.wikipedia.org/wiki/Receiver_operating_characteristic).
- [147] Wikipedia. Rejection sampling. [https://en.wikipedia.org/wiki/Rejection\\_sampling](https://en.wikipedia.org/wiki/Rejection_sampling).
- [148] Wikipedia. Score test. [https://en.wikipedia.org/wiki/Score\\_test](https://en.wikipedia.org/wiki/Score_test).
- [149] Wikipedia. Signal flow graph. [https://en.wikipedia.org/wiki/Signal-flow\\_graph](https://en.wikipedia.org/wiki/Signal-flow_graph).
- [150] Wikipedia. Simple linear regression. [https://en.wikipedia.org/wiki/Simple\\_linear\\_regression](https://en.wikipedia.org/wiki/Simple_linear_regression).
- [151] Wikipedia. Simpson's paradox. [https://en.wikipedia.org/wiki/Simpson's\\_paradox](https://en.wikipedia.org/wiki/Simpson's_paradox).
- [152] Wikipedia. Spring system. [https://en.wikipedia.org/wiki/Spring\\_system](https://en.wikipedia.org/wiki/Spring_system).
- [153] Wikipedia. Student's t-distribution. [https://en.wikipedia.org/wiki/Student's\\_t-distribution](https://en.wikipedia.org/wiki/Student's_t-distribution).
- [154] Wikipedia. Support vector machine. [https://en.wikipedia.org/wiki/Support-vector\\_machine](https://en.wikipedia.org/wiki/Support-vector_machine).
- [155] Wikipedia. Survival analysis. [https://en.wikipedia.org/wiki/Survival\\_analysis](https://en.wikipedia.org/wiki/Survival_analysis).

- [156] Wikipedia. Time series. [https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series).
- [157] Wikipedia. Transfer learning. [https://en.wikipedia.org/wiki/Transfer\\_learning](https://en.wikipedia.org/wiki/Transfer_learning).
- [158] Wikipedia. Transformer (machine learning model). [https://en.wikipedia.org/wiki/Transformer\\_\(machine\\_learning\\_model\)](https://en.wikipedia.org/wiki/Transformer_(machine_learning_model)).
- [159] Wikipedia. Uplift modelling. [https://en.wikipedia.org/wiki/Uplift\\_modelling](https://en.wikipedia.org/wiki/Uplift_modelling).
- [160] Wikipedia. Variational Bayesian methods. [https://en.wikipedia.org/wiki/Variational\\_Bayesian\\_methods](https://en.wikipedia.org/wiki/Variational_Bayesian_methods).
- [161] Wikipedia. Vector autoregression. [https://en.wikipedia.org/wiki/Vector\\_autoregression](https://en.wikipedia.org/wiki/Vector_autoregression).
- [162] Wikipedia. Viterbi algorithm. [https://en.wikipedia.org/wiki/Viterbi\\_algorithm](https://en.wikipedia.org/wiki/Viterbi_algorithm).
- [163] Wikipedia. Wald test. [https://en.wikipedia.org/wiki/Wald\\_test](https://en.wikipedia.org/wiki/Wald_test).
- [164] Wikipedia. Z-transform. <https://en.wikipedia.org/wiki/Z-transform>.
- [165] Hao Wu and Zhaohui Steve Qin. course notes, BIOS731: Advanced statistical computing, 2016 Emory Univ. <http://web1.sph.emory.edu/users/hwu30/teaching/statcomp/statcomp.html>.
- [166] Ronghui (Lily) Xu. Lecture notes, MATH 284, Spring 2020, Survival analysis. <https://mathweb.ucsd.edu/~rxu/math284/>.
- [167] Jonathan S Yedidia, William T Freeman, and Yair Weiss. Understanding belief propagation and its generalizations, Mitsubishi Technical Report tr-2001-22. <https://merl.com/publications/docs/TR2001-22.pdf>.