# Index

0-1 loss, **104**, 276

Absolute value rectification, 192
Accuracy, 425
Activation function, 170
Active constraint, 95
AdaGrad, 307
ADALINE, *see* adaptive linear element
Adam, 308, 427
Adaptive linear element, 15, 24, 27
Adversarial example, 268
Adversarial training, 268, 271, 532
Affine, 110
AIS, *see* annealed importance sampling
Almost everywhere, 71
Almost sure convergence, 130
Ancestral sampling, 582, 597
ANN, *see* Artificial neural network
Annealed importance sampling, 627, 670, 719
Approximate Bayesian computation, 718
Approximate inference, 585
Artificial intelligence, 1
Artificial neural network, *see* Neural network
ASR, *see* automatic speech recognition
Asymptotically unbiased, 124
Audio, 102, 360, 460
Autoencoder, 4, 356, **504**
Automatic speech recognition, 460

Back-propagation, 203
Back-propagation through time, **384**
Backprop, *see* back-propagation

Bag of words, 473
Bagging, 256
Batch normalization, 268, 427
Bayes error, **117**
Bayes' rule, 70
Bayesian hyperparameter optimization, 438
Bayesian network, *see* directed graphical model
Bayesian probability, 55
Bayesian statistics, **135**
Belief network, *see* directed graphical model
Bernoulli distribution, 62
BFGS, 316
Bias, 124, 229
Bias parameter, 110
Biased importance sampling, 595
Bigram, 464
Binary relation, 484
Block Gibbs sampling, 601
Boltzmann distribution, 572
Boltzmann machine, 572, 656
BPTT, *see* back-propagation through time
Broadcasting, 34
Burn-in, 599

CAE, *see* contractive autoencoder
Calculus of variations, 179
Categorical distribution, *see* multinoulli distribution
CD, *see* contrastive divergence
Centering trick (DBM), 675
Central limit theorem, 63
Chain rule (calculus), 206
Chain rule of probability, 59