# INDEX

weight
    naming convention for
      neurons, 322
    for neurons, 316
    overview in deep network, 11
weight sharing, 433
weighted coin, 86
weighted graph, 324
weighted plurality voting, 299
weirdness (high-dimensional), 175
width (of a recurrent cell), 551
winner (egg), 160
wire (in a neural network), 323
word embedding, 566

## X

Xavier Normal initialization, 325
Xavier Uniform initialization, 325

## Y

yolker, 160
Yorkshire terrier, 184

## Z

zero gradient, 128
zero padding, 442
zero point, 436
zero-dimensional array, 328
zero-shot training, 594