

Supervised learning: labeled with correct behavior

Reinforcement: learning interacts with the world, max scalar reward

Unsupervised: no label example - looking for interesting pattern

KNN: $\|x^{(a)} - x^{(b)}\|_2 = \sqrt{\sum_{j=1}^d (x_j^{(a)} - x_j^{(b)})^2}$ j -dimension $x \in \mathbb{R}^d$

Algo: 1. Find k examples $\{x^{(i)}, t^{(i)}\}_{i=1}^N$ closest to the test instance x .
2. Classification output is major class.

$y^* = \max_{t^{(i)} \in \text{class label}} \sum_{i=1}^k I(t^{(i)} = t^{(i)})$ data

small k: may overfit, be sensitive to random idiosyncrasies in training
large k: may underfit, fail to capture important regularities.
balance: depends on $n, k < \sqrt{n}$, choose use validation set.

measure generalization error using test set. each point is $O(d)$

Curse of dimensionality: The volume of a single ball of radius ϵ around x is $\pi d \epsilon^d$. The total volume of $[0, 1]^d$ is 1 . $O(\frac{1}{\epsilon})^d$ points are needed to cover the volume. In high dimension, most points are approx. same distance.

Normalization: sensitive to different ranges of different features (units)

$\tilde{x}_j = \frac{x_j - \mu_j}{\sigma_j}$ simple fix: normalize each dimension to be zero mean & unit variance.

Computation cost: compute D -dimension Euclidean dis with N^2 points $O(DN^2)$
Sort dis $O(N \log N)$ · Store entire dataset · Done for each query

Entropy: High Entropy · flat hist, uncertainty ↑, unpredictable

$P(y|x) = \frac{P(x,y)}{P(x)}$ $P(x) = \sum_y P(x,y)$ $H(Y) = -\sum_{y \in Y} P(y) \log_2 P(y)$

Joint Entropy: $H(X,Y) = -\sum_{x \in X} \sum_{y \in Y} P(x,y) \log_2 P(x,y)$

Conditional Entropy: $H(Y|X=x) = -\sum_{y \in Y} P(y|x) \log_2 P(y|x)$

Expected: $H(Y|X) = \sum_{x \in X} P(x) H(Y|X=x)$

$= -\sum_{x \in X} \sum_{y \in Y} P(x,y) \log_2 P(y|x)$

$H(X,Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$

X, Y independent: $H(Y|X) = H(Y)$ $H(Y|Y) = 0$ $H(Y|X) \leq H(Y)$

$IG(Y|X) = H(Y) - H(Y|X)$ X completely uninformative $IG(Y|X) = 0$ informative $IG(Y|X) = H(Y)$

Algo: Decision Tree
Pick a feature to split at a non terminal node
Split examples into groups based on feature value
for each group: if no example, return majority from parents
else if all example in same class - return class
else loop to step 1
Terminal when all leaves contain same class or empty.

{ not too small: need to handle important but possible subtle distinction in data ; not too big: computation efficiency, avoid overfitting }
Human Interpretability \Rightarrow desire small tree with info nodes near root.

KNN { fewer hyperparameter (k) ; incorporate interesting distance measure DT } simple to deal with discrete features, missing values, and poorly scaled data. Fast at test time. more interpretable

Bias-Variance Decomposition:

Generalization: Ability of a model to correctly classify/predict from unseen example (from the same dist, that the training data was drawn)

Square Error loss: $L(y, t) = \frac{1}{2}(y - t)^2$ $P(t|x)$: random, choose y

$y^* = E[t|x]$ prove: $E[(y-t)^2|x] = E[y^2 - 2yt + t^2|x] = y^2 - 2yE[t|x] + E[t^2|x]$

$\frac{\partial J(w)}{\partial w} = \frac{1}{2} [2x^T(\bar{t} - x\bar{w})] = x^T x \bar{w} - x^T \bar{t}$ set $\frac{\partial J}{\partial w} = 0 \Rightarrow \bar{w} = (x^T x)^{-1} x^T \bar{t}$

$= y^2 - 2yE[t|x] + E[t|x]^2 + \text{Var}[t|x] \quad \text{Var}[x] = E[x^2] - E[x]^2$

$= y^2 - 2y\bar{y} + \bar{y}^2 + \text{Var}[t|x] = (y - \bar{y})^2 + \text{Var}[t|x]$ Bayes Error

Bayes Optimal: Best we hope to with any learning algo.
Treat y as random: $E[y-t] = E[(y-\bar{y})^2] + \text{Var}[t]$

$= E[y^2 - 2y\bar{y} + \bar{y}^2] + \text{Var}[t] = (\bar{y}^2 - 2\bar{y}E[y] + E[y^2]) + \text{Var}[y] + \text{Var}[t]$

Bayes Optimality: $E[y-t] = (y^* - E[y])^2 + \text{Var}[y] + \text{Var}[t]$ (batch)

Bias: How wrong the \hat{y} prediction is corresponds to underfitting.

Variance: the amount of variability in predictions (overfitting)

Bayes error: the inherent unpredictability of the targets.

prediction y at query x is a random variable (from the choice of dataset)

y^* is the optimal deterministic prediction (choice of label)

t is a random target from sample from true conditional $P(t|x)$ (from dataset)

Bagging: sample m independent training set from Psample

Idea { compute y_i using each training set
compute average prediction $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$

Bias unchanged $E[\bar{y}] = E[\frac{1}{m} \sum_{i=1}^m y_i] = \frac{1}{m} \cdot m \cdot E[y_i] = E[y_i]$

Variance ↓: $\text{Var}[\bar{y}] = \text{Var}\left[\frac{1}{m} \sum_{i=1}^m y_i\right] = \frac{1}{m^2} \cdot m \cdot \text{Var}[y_i] = \frac{1}{m} \text{Var}[y_i]$

Algo: Given training set D , use the empirical distribution P_D as a proxy for Psample ("bootstrap aggregation or bagging")

Take a dataset D with n examples · Generate m new dataset ("resample or bootstrap sample"). Each dataset has n examples sampled from D with replacement · Average the predictions of models trained on the m dataset. As $|D| \rightarrow \infty$, $P_D \rightarrow$ Psample / foreach node, choose a random subset of features

Random forest: reduce correlation between bagged decision tree. Var ↓

Bagging: reduce overfitting by averaging predictions · A small ensemble often better than a single great model. limitation: · bias unchanged in case of square error. · Correlation in classifiers mean less $\text{Var} \downarrow \Rightarrow$ Add more randomness in RF. · Weighted members equally may not be the best \Rightarrow weighted ensambing often leads to better results if member are very different. ($y_{\text{bagged}} = \frac{1}{m} \sum_{i=1}^m y_i > 0.5$) Binary

Linear regression: $y = f(x) = \sum_{j=1}^D w_j x_j + b = \bar{w}^T x + b$ (w, b are parameter)

Loss function: $L(y, t)$ defines how badly the algo's prediction y fits target t for same example x . Square error loss func: $L(y, t) = \frac{1}{2}(y - t)^2$ residual $y - t$

Cost function: loss function averaged over all training example L · empirical or average "loss".

$J(w, b) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2 = \frac{1}{2N} \sum_{i=1}^N (w^T x^{(i)} + b - t^{(i)})^2$ dataset.

Design matrix & target vector \Rightarrow compute prediction for the whole dataset.

$x^T w + b = y$ $(x_1^{(1)} \ x_2^{(1)} \ \dots \ x_D^{(1)}) \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{pmatrix} + b \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = (y^{(1)} \ y^{(2)} \ \dots \ y^{(N)})$ vs $\begin{pmatrix} t^{(1)} \\ t^{(2)} \\ \vdots \\ t^{(N)} \end{pmatrix}$

$y = x^T w \quad x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_D^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D+1}$ and $w = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \in \mathbb{R}^{D+1}$

Squared error cost: $\bar{y} = x^T w + b$, $J = \frac{1}{2N} \|y - \bar{y}\|^2$, $J = \frac{1}{2} \|y - \bar{y}\|^2$

Optimization: minimize cost function $J(\bar{w})$

$\nabla_{\bar{w}} J = \frac{\partial J}{\partial w} = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_D} \end{pmatrix}$

① direct: Take gradient wrt w , set it = 0
② iterative solution sum:

$\nabla_{\bar{w}} J(w) = \frac{1}{2} \|x^T w - \bar{t}\|^2 = \frac{1}{2} (\bar{t} - x\bar{w})^T (\bar{t} - x\bar{w}) = \frac{1}{2} (\bar{t} - x\bar{w})^T \frac{1}{2} (\bar{t} - x\bar{w})$

$\frac{\partial J(w)}{\partial w} = \frac{1}{2} [2x^T(\bar{t} - x\bar{w})] = x^T x \bar{w} - x^T \bar{t}$ set $\frac{\partial J}{\partial w} = 0 \Rightarrow \bar{w} = (x^T x)^{-1} x^T \bar{t}$

③ Initialize the weights to something reasonable (e.g. all zeros) and repeatedly adjust them in the direction of steepest descent.

$\frac{\partial J}{\partial w_j} > 0 \Rightarrow$ decreasing J requires decreasing w_j $\frac{w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j}}{\partial w_j}$

$\frac{\partial J}{\partial w_j} < 0 \Rightarrow$ decreasing J requires increasing w_j $\frac{w_j \leftarrow w_j + \alpha \frac{\partial J}{\partial w_j}}{\partial w_j}$

learning rate: α ($d \uparrow$, change fast) If min total loss $\Rightarrow \alpha^* = \alpha/N$ $\frac{0.01 \sim 0.1}{\alpha}$

$\bar{w} \leftarrow \bar{w} - \alpha \frac{\partial J}{\partial \bar{w}}$ (for: Linear Re: $\bar{w} \leftarrow \bar{w} - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x^{(i)}$)

Gradient descent: update \bar{w} in the direction of fastest decrease.

avg: $\frac{\partial J(w)}{\partial w} = \frac{1}{N} x^T (x\bar{w} - \bar{t})$ $\begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D}$

$J(w) = \frac{1}{2} \|x\bar{w} - \bar{t}\|^2 = \frac{1}{2} \sum_{i=1}^N x^{(i)} \cdot (y^{(i)} - t^{(i)})$ $\geq \sum_{i=1}^N x^{(i)} \cdot [y^{(i)} - t^{(i)}]$

(sum of $d/1$ d dimensional input vectors) $\boxed{[d \times 1]}$

Why? · more efficient in high dimensional space, matrix inversion $O(D^3)$
· Gradient descent update $O(ND)$ or less with stochastic ·
· easier to implement · applicable to a much broader set of models

Feature mapping: In model non linear relationship map input features to another space $\Psi(x) : \mathbb{R}^D \rightarrow \mathbb{R}^d$, linear

Polynomial Feature mapping: $y = w_0 + w_1 x + w_2 x^2 + \dots + w_m x^m = \sum_{i=0}^m w_i x^i$

Feature mapping: $\Psi(x) = [1, x, x^2, \dots, x^m]^T$, $y = \Psi^T(x) \bar{w}$ is linear in \bar{w} .

M -too small - too simple - underfit · too large - too complex - overfit

Regularization: control the model complexity, improve generalization L^2 penalty: as our regularizer $R(w) = \frac{1}{2} \|\bar{w}\|_2^2 = \frac{1}{2} \sum w_j^2$ fit data

$J_{\text{reg}}(\bar{w}) = J(\bar{w}) + \lambda R(\bar{w}) = J(\bar{w}) + \frac{\lambda}{2} \sum w_j^2$ tradeoff \leq norm of ws

$\bar{w}_{\text{Ridge}} = \arg \min_{\bar{w}} J_{\text{reg}}(\bar{w}) = \arg \min_{\bar{w}} \frac{1}{2N} \|\bar{w}\|_2^2 + \frac{\lambda}{2} \|\bar{w}\|_2^2 = (x^T x + \lambda N I)^{-1} x^T \bar{t}$

weight decay: $\bar{w} \leftarrow \bar{w} - \alpha \frac{\partial}{\partial \bar{w}} (J + \lambda R) = \bar{w} - \alpha (\frac{\partial J}{\partial \bar{w}} + \lambda \bar{w})$

concave \downarrow $\bar{w} - \alpha (\frac{\partial J}{\partial \bar{w}} + \lambda \bar{w})$
convex \uparrow $\bar{w} - \alpha (\frac{\partial J}{\partial \bar{w}} + \lambda \bar{w}) = (1 - \alpha \lambda) \bar{w} - \alpha \frac{\partial J}{\partial \bar{w}}$

Convex: if a function is convex, then every critical point is global optimum

use set S is convex if any line segment connecting two points in S lies entirely within S

$x_1, x_2 \in S \Rightarrow \lambda x_1 + (1-\lambda)x_2 \in S$ for $0 \leq \lambda \leq 1 \Rightarrow \lambda x_1 + \dots + \lambda N x_N \in S$ ($\lambda_1 + \dots + \lambda_N = 1$)

weighted averages, convex combinations of points in S lies within $x_1, \dots, x_N \in S$

convex function: $f((1-\lambda)x_0 + \lambda x_1) \leq (1-\lambda)f(x_0) + \lambda f(x_1)$ bowl shape Δ

$Z = w^T x + b$, if loss func is a convex func of $Z \Rightarrow$ convex of w, b

J · large, overshoot, oscillations, missing optimum, instability
· small α , slow convergence
choose α : · Grid search
good convergence
iteration → validation set · Training curve

Stochastic Gradient descent (batch: sum over entire dataset)
updates the parameters based on the gradient for one training example

Repeat { (1) choose example i uniformly at random
(2) perform update $\theta \leftarrow \theta - \alpha \frac{\partial L^{(i)}}{\partial \theta}$

Pros: · cost of each update is independent of N · make sign progress before seeing all the data · stochastic gradient descent is an unbiased estimate of the batch gradient given sampling each example uniformly at random

$\frac{\partial J(w)}{\partial w} = \text{Ex} \left[\frac{\partial L^{(i)}}{\partial w} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial L^{(i)}}{\partial w}$ · High variance, can't exploit efficient vectorized operations

Mini-Batch Gradient descent: randomly chosen medium-sized subset of training examples M . In theory, sample examples independently and uniformly with replacement. In practice, permute the training set and then go through it sequentially. Each path over the data is called an epoch.

{ large size: more computation time, estimate accurate, exploit vectorization
small size: faster update, estimate noisier, cannot exploit vectorization
 M -hyperparameter, $M=100$ reasonable. α - influence noise in parameter
start with large α to get closer to optimal, gradually decrease the learning rate to fluctuation }

Binary Linear Classification $Z = W^T X + b$ $y = \begin{cases} 1 & \text{if } Z \geq r \\ 0 & \text{if } Z < r \end{cases}$

Simplified: $Z = W^T X$, $y = \begin{cases} 1 & \text{if } Z \geq 0 \\ 0 & \text{if } Z < 0 \end{cases}$ ($W^T X + b \geq 0$ add $x_0 = 1$, $w_0 = b$)

NOT: $x_0 = 1$, $x_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $x_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, $x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $x_4 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}$, $x_5 = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$

The line $\{x : W^T x = 0\}$ defines the **decision boundary**. A line in 2D, or a hyperplane in high dimensions. pass through origin.

Data is linearly separable if a linear decision rule can perfectly separate the training examples. In to guarantee correct classification weight space: each training example specifies a half space that w must lie feasible region satisfy all the constraints.

The problem is **feasible** if the feasible region is non empty.

Data space: points are training examples, a line represents a set of weights. Weight space: points are weights, a half space represents a training example constraining the weights. Linear programming

Learn weights for linearly separable dataset $\{0\}$ perceptron algorithm

Logistic Regression: if data not linearly separable • define loss func. find weights that minimize the average loss over the training examples.

0-1 Loss: $L_{0-1}(y, t) = \begin{cases} 0, & \text{if } y=t \\ 1, & \text{if } y \neq t \end{cases}$ cost $J = \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, t^{(i)})$ cons: to mini loss, we need critical point, but ∇ is zero almost everywhere $\Rightarrow \Delta w$ has no effect on the loss. 0-1 Loss is discontinuous at $Z=0$, ∇ is undefined.

squared loss: $LSE(Z, t) = \frac{1}{2}(Z - t)^2$. Treat binary targets as continuous values. Make final prediction y by thresholding Z at $\frac{1}{2}$.

cons: if $t=1$, $L(Z=1) > L(Z=0) \Rightarrow$ Making a correct prediction with high confidence should be good. but incur a large loss.

Logistic Activation function: squash prediction y into $[0, 1]$, $y = \sigma(z) = \frac{1}{1+e^{-z}}$ (activation func.)

$Z = W^T X$, $y = \sigma(Z)$, $LSE(y, t) = \frac{1}{2}(y - t)^2$, $\sigma'(z) = \frac{1}{(1+e^{-z})^2}$ (activation func.)

cons: suppose $t=1$ and Z is very negative $\Rightarrow y \approx 0$ is very wrong, but the weights appears to be at the critical point: $\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_j} = 0 \cdot \frac{\partial z}{\partial w_j} = 0$

Cross Entropy loss: interpret $y \in [0, 1]$ as the prob. $t=1$. Heavily penalize the extreme mis-classification cases. $LCE = -t \log y - (1-t) \log(1-y)$

Numerical Instability: $t=1, Z \neq 0, y \rightarrow 0 \log y \rightarrow \log 0$ assume $t=1$

Logistic Cross Entropy: $LLCE(Z, t) = t \log(1+e^{-Z}) + (1-t) \log(1+e^{-Z})$

$E = f^* + np.\logaddexp(0, -Z) + (1-t) * np.\logaddexp(0, Z)$

Gradient of Logistic loss: $\frac{\partial LCE}{\partial w_j} = \frac{\partial LCE}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_j} = (-\frac{1}{y} + \frac{1-t}{1-y}) y(1-y) x_j = (y - t)y x_j$

$w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j} = w_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$ generalized linear model

Gradient Checking with finite difference: $\frac{\partial}{\partial x_i} f(x_1, \dots, x_N) = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_i+h, \dots, x_N) - f(x_1, x_2, \dots, x_i-h, \dots, x_N)}{2h}$

eg. $x_1=3, x_2=7, h=0.01 = \frac{\partial f}{\partial x_1} \approx f(3, 7.01) - f(3, 6.99)$

Run gradient checks on small, randomly chosen inputs.

use double precision floats • plug in small value of h as 10^{-10}

compute relative error = $\frac{|a-b|}{|a+b|}$ a : finite difference b : derivative by your implement (should be very small 10^{-6})

Hypothesis: is a function $f: X \rightarrow T$ from the input to target space

Hypothesis space H as a set of hypothesis (eg a set of linear functions)

Inductive bias: preference for some h of H over others better than others

No free lunch: if dataset/problem were not naturally biased \Rightarrow no ML would be

parametric algo = H is defined using a finite set of parameters eg DT

non parametric = hypothesis H is defined in terms of the data eg KNN.

is linearity, constancy no parameters defined in terms of data.

Multi-class classification: Targets form a discrete set $\{1, \dots, k\}$

represent targets as one-hot vectors or one of k encoding. $\vec{t} = (t_0, \dots, t_k) \in \mathbb{R}^k$ don't add explicit $R(D)$ to cost

vectorized form: $\vec{Z} = W\vec{x} + \vec{b}$ or $Z = W\vec{x}$ (with $x_0 = 1$)

Non vectorized form: $Z_k = \sum_{j=1}^k w_{kj} x_j + b_k$ for $k = 1, 2, \dots, k$

$W: k \times D$, $X: D \times 1$, $b: k \times 1$, $Z: k \times 1$ the belief that model has class k

$D=2, k=4$, $W = \begin{bmatrix} -w_1 \\ -w_2 \\ -w_3 \\ -w_4 \end{bmatrix}, \vec{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}, \vec{Z} = W\vec{x} + \vec{b} = \begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ Z_4 \end{bmatrix}$

(Z_k): how much the model prefers the k th prediction. $W^T x + b_k$

$y_i = \begin{cases} 1, & \text{if } i = \arg \max_k Z_k \\ 0, & \text{otherwise} \end{cases}$

Softmax Regression: $y_k = \text{softmax}(Z_1, \dots, Z_k)_k = \frac{e^{Z_k}}{\sum_k e^{Z_k}}$

$\sum_k y_k = 1$ • Z_k are called logits • y_k as probability • If Z_k is much larger than the others, then $\text{softmax}(Z_k)_k \approx 1$ and it behaves like $\arg \max$

Cross Entropy Loss: $LCE(y, t) = -\sum_{k=1}^K t_k \log y_k = -\vec{t}^T (\log \vec{y})$ to loss.

$t^T = [0 \ 1 \ 0 \ 0]$, $\vec{y}^T = [0.01 \ 0.9 \ 0.05 \ 0.04]$ only one output contribute

Gradient descent for softmax: $\frac{\partial LCE}{\partial w_k} = \frac{\partial LCE}{\partial Z_k} \cdot \frac{\partial Z_k}{\partial w_k} = (y_k - t_k) \vec{x}$

$w_k \leftarrow w_k - \alpha \cdot \frac{1}{N} \cdot \sum_{i=1}^N (y_k^{(i)} - t_k^{(i)}) \vec{x}^{(i)}$ other matrix

model performance: • **Progress** = training curve ($\#$ of iteration vs error) • **Accuracy** (average 0-1 loss, error rate) $\text{Acc} = \frac{TP + TN}{P+N}$ true class pred. o TN FPI TP II highly sensitive to class imbalance = class I FPI

sensitivity = $TP / (TP + FN)$ **specificity** = $TN / (TN + FP)$

Receiving Operating Characteristic (ROC) Curve: TP frequency

Area under curve can quantify if a binary classifier achieves a good tradeoff between S, S

Confusion matrix: rows are true label, cols are predicted labels, entries

Limits of linear classification: Some dataset are not linearly separable. (Prove Separable): Each training point \rightarrow linear rule \rightarrow inequality constraint in W space \rightarrow feasible region. (Prove not separable): Proof by contradiction

Half space are convex: if two points lie in the same space, the line segment connecting them also lie in the same space/half. Suppose the problem is feasible if positive examples are in the positive half-space, then the green line segment must be as well. Similarly, the red line must in the negative half-space But the intersection can't lie in both half-spaces. Contradiction!

Feature map: $\phi(x) = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix} \mid \phi_1(x_1), \phi_2(x_2), \dots, \phi_D(x_D) \mid + |$

Neural Network: multi-layer perceptrons consists fully connected layer

Fully connected layer: all input units are connected to all output units

Every layer i : weight matrix $N_i \times N_{i-1}$, connect N_{i-1} input N_i output units. $y = \phi(WX + b) \Rightarrow \phi$: applied component wise (non-linearity)

$y = \phi(x) = \phi(WX + b)$

Activation function: Identity, Rectified Linear Unit (ReLU), Hyperbolic tangent, Soft ReLU, Logistic, Hard threshold

Compute of each layer: last layer \leftarrow regression $y = f^{(L)}(h^{(L-1)}) = (W^{(L)})^T h^{(L-1)} + b^{(L)}$ binary $y = f^{(L)}(h^{(L-1)}) = \sigma((W^{(L)})^T h^{(L-1)} + b^{(L)})$

Expressivity: two model A, B, $A \gg B$ \Rightarrow A is more expressive can represent any fun in H_B

Note: Linear NN $y = W^{(3)} W^{(2)} W^{(1)} \vec{x}$ can only represent linear functions, no more expressive than linear regression.

Non Linear NN with non linear activation functions. Universal func

Approximator: can approximate any function arbitrarily well For any $f: X \rightarrow T$, there is a sequence of $f_i \in H$ with $f_i \rightarrow f$.

Binary Input & Targets: 2^D hidden units, each of which responds to one particular input configuration, require only 1 hidden layer

Early stopping to avoid overfitting Error/Training error # iteration. $\frac{1}{2} \frac{\partial^2 E}{\partial \theta^2}$

Error Signal: $\vec{y} = dV/dy$, just \vec{a} values our prog computing

Multivariate chain rule: $\frac{d}{dt} f(X(t), Y(t)) = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$ parent child.

Back Propagation: $\vec{t} = \vec{x} \frac{\partial x}{\partial t} + \vec{y} \frac{\partial y}{\partial t}$ parents before child

let V_1, \dots, V_N be a topological ordering of the computation graph. node, we try to compute gradients.

Forward pass: For $i=1, \dots, N$, compute V_i as a function of Parent(V_i)

Backward Pass: For $i=N-1, \dots, 1$, $\vec{V}_i = \sum_j \vec{V}_j \frac{\partial V_i}{\partial V_j}$

Regularized logistic least squares:

Forward: $Z = W^T X + b$, $y = \sigma(Z)$, $L = \frac{1}{2} (y - t)^2$, $R = \frac{1}{2} \|W\|^2$, $L_{reg} = L + \lambda R$

Backward: $\vec{y} = \vec{L} \cdot \frac{dL}{dy} = \vec{L}(y - t)$, $\vec{z} = \vec{y} \frac{dy}{dz} = \vec{y} \sigma'(Z)$, $\vec{r} = \vec{L}_{reg} \frac{dL_{reg}}{dr} = \vec{L}_{reg} \frac{dL}{dr}$, $\vec{w} = \vec{z} \frac{dz}{dw} + \vec{r} \frac{dr}{dw}$, $\vec{b} = \vec{z} \frac{dz}{db} = \vec{z}$, $\vec{L}_{reg} = \vec{L} + \lambda \vec{R}$

BP for two layer NN:

Forward: $Z_i = \sum_j W_{ij}^{(1)} x_j + b_i^{(1)}$, $h_i = \sigma(Z_i)$, $y_k = \sum_i W_{ki}^{(2)} h_i + b_k^{(2)}$, $\vec{w}_{ij}^{(1)} = \vec{w}_{ki}^{(2)}$, $\vec{b}_i^{(1)} = \vec{b}_k^{(2)}$

Back: $\vec{z}_i = \sum_k \vec{y}_k W_{ki}^{(2)}$, $\vec{h}_i = \vec{z}_i \sigma'(Z_i)$, $\vec{w}_{ij}^{(2)} = \vec{z}_i \vec{w}_{ij}^{(1)}$, $\vec{b}_i^{(2)} = \vec{b}_i^{(1)}$

Forward: $Z_i = \sum_j W_{ij}^{(1)} x_j + b_i^{(1)}$, $h_i = \sigma(Z_i)$, $y_k = \sum_i W_{ki}^{(2)} h_i + b_k^{(2)}$, $\vec{w}_{ij}^{(1)} = \vec{w}_{ij}^{(2)}$, $\vec{b}_i^{(1)} = \vec{b}_i^{(2)}$

Back: $\vec{z}_i = \sum_k \vec{y}_k W_{ki}^{(2)}$, $\vec{h}_i = \vec{z}_i \sigma'(Z_i)$, $\vec{w}_{ij}^{(2)} = \vec{z}_i \vec{w}_{ij}^{(1)}$, $\vec{b}_i^{(2)} = \vec{b}_i^{(1)}$

Forward: one add-multiply op per weight

Back: two add-multiply op per weight

Forward: one add-multiply op per weight

Back: quadratic in units per layer

Forward: linear in # of layer

Back: quadratic in units per layer

vectorized: $W^{(1)} \otimes x \rightarrow W^{(2)} \rightarrow t \rightarrow L = \frac{1}{2} \|t - y\|^2$, $\vec{w}^{(1)} = \vec{y} h^T$, $\vec{w}^{(2)} = \vec{h} \vec{w}^{(1)}$, $\vec{b}^{(2)} = \vec{y}$, $\vec{w}^{(1)} = \vec{z} X^T$, $\vec{b}^{(1)} = \vec{z}$

Fully connected layer: all input units are connected to all output units

Every layer i : weight matrix $N_i \times N_{i-1}$, connect N_{i-1} input N_i output units. initial, $y = W^{(2)} h + b^{(2)}$

should not set all $W=0$; $L = \frac{1}{2} \|t - y\|^2$, neural net implausible.

Invariant: image transformed slightly \Rightarrow classification unchanged

convolution ID: $(X * W)[t] = \sum_{s,t} X[s, t] W[t-s, t-t]$

commutativity • linearly $\vec{w}^{(1)} = \vec{z} X^T$, $\vec{w}^{(2)} = \vec{h} \vec{w}^{(1)}$

2D: $(X * W)[i, j] = \sum_s \sum_t X[i-s, j-t] * W[s, t]$

max pooling: $y_i = \max_j z_j$ in pooling group reduce size

$(AB)^T = A^T B^T$, $(AB)^{-1} = B^{-1} A^{-1}$, $\frac{\partial X^T A}{\partial x} = \frac{\partial A^T}{\partial x} = A$, $\frac{\partial X^T B X}{\partial x} = (B + B^T) X$, $\frac{\partial (X - AS)^T W (X - AS)}{\partial x} = -2A^T W (X - AS)$

$(RA)^T = R A^T$, $I^{-1} = I$, $(A B)^T = B^T A^T$, $I^T = I$, $f(W) = W^T A W + b^T W + y$, $\nabla f = \frac{1}{2} (A^T + A) W + b$, $AA^T = A^T A = I$, $(RA)^{-1} = I - A^T R^{-1}$, $\frac{\partial f}{\partial x} = AW + b$ (if A sym), $\frac{\partial f}{\partial x} = 2W(X - AS)$, $\frac{\partial^2 f}{\partial x^2} = A + A^T = 2A$, $\frac{\partial f}{\partial x} = -2W(X - AS)^T$