```
1 Review
2 ---
3
4 1. Perceptron(Linear separator)
5     - Hypothesis set: sign(W^Tx)
6     - Learning algorithm: PLA
7         - idea: start with some weights and try to
   improve it
8         - algorithm:
9         ```
10         for i in max_iter:
11             for (x, y) in in_samples:
12                 if sign(W^Tx) != y:
13                     all_mispred_samples.append((x,
   y))
14             for (x, y') in all_mispred_samples:
15                 W_(i-1) += y * x
16         ```
17         - pocket algorithm: use the model with the
   smallest in-sample error
18     - Theoretical basis: If the data can be fit by
   a linear separator, then after some finite number
   of steps, PLA will find one.
19
20 2. Linear Regression
21     - y = W^Tx, E = (y-W^Tx)^2
22     - Has analytic solution:
23         - Normal equation: take the derivative of E
   , we get X^TXw = X^Ty
24         - The linear regression algorithm gets the
   smallest possible Ein in one step
25
26 3. Logistic Regression(for classification task)
27     - predict the probability with: sigmoid(W^Tx)
28     - sigmoid: 1 / (1 + e^(-s))
29     - Cross entropy loss(for binary case): E = ln(1
   + e(-y * W^T * x))
30         - y takes {+1, -1}, y is the real label
31     - Cross entropy is an alternative
   representation of maximum likelihood estimation,
```

```
31 the loss function used in logReg
32     is actually derived from MLE
33
34 4. Softmax Regression(for multi-class
   classification)
35     - formula:
36     - relationship with logistic regression:
37         - equivalence between softmax and sigmoid
38         - equivalence between binary cross entropy
   and multi-class cross entropy:
39     - Cross entropy loss(for multi-class case)
40
41 5. Gradient Descent
42     - Batch-GD: update for the entire batch
43     - SGD: update for each sample
44     - Mini-Batch GD: update for each mini-batch (
   accumulate loss -> compute gradient -> update)
45
46 6. Neural Network
47     - Backpropagation Algorithm
48     ```
49     Compute_X(sample):
50         s, x = new_array(), new_array()
51         s[0] = sample
52         for l in (1, L):
53             x[l] = s[l-1] * W[l]
54             s[l] = actv(x[l])
55         return x
56     ```
57     ```
58     Compute_sensitivity(x):
59         E' = Loss'(S[L])
60         sen[L] = E' * actv'(s[L])
61         for l in (L-1, 1):
62             sen[l] = sen[l+1] * W[l+1]^T * actv'(x[
   l])
63         return sen
64     ```
65     ```
66     Backprop():
```

```
67     for s in all_samples:
68         x = Compute_X(s)              //
   compute output of each layer
69         sen = Compute_sensitivity(s) //
   compute sensitivity of each layer
70         E += Loss(s)
71         for l in L:                   // for
   each layer
72             G[l] = x[l-1] * sen[l-1]  //
   compute gradient for layer l
73             G[l] += G[l] + 1/N * G[l]  //
   accumulate gradient
74     for l in L:
75         W[l] -= lr * G[l]              //
   update weight
76     ```
77     - Generalization
78         - L1 regularization
79         - L2 regularization (weight decay)
80         - Early stopping
81         - Dropout
82         - Data augmentation
83     - Better GD
84         - variable learning rate
85         - Steepest Descent (Line Search): binary
   search to decide the lr that minimize E in one
   step
86         - ... and others
87
88 7. CNN
89     - Domain knowledge
90         - translation invariance
91         - locality
92     - Basics (hk is the size of the kernel)
93         - Conv
94             - H(i+1) = H(i) + hk -1
95             - Backprop(DeConv): 1. full-padding 2
   . conv with inverted filters
96         - Padding
97             - "same" padding (or half-padding): Ph
```

```
97 = (hk - 1)/2 on each side
98             - "valid" padding: not use any padding
   ...
99             - without padding, pixels on
   boarders are under-represented
100         - "full" padding: Ph = hk - 1 on each
   side, increase by hw-1, useful when doing reverse
   convolution
101         - Strides (Sk is the stride)
102             - H(i+1) = (H(i) - hk) / Sk + 1
103         - Bias
104             - one for each channel, added on each
   pixel
105         - Pooling
106             - diff with conv: apply on each
   feature map, does not change the number of feature
   maps
107             - with stride of 1, H(i+1) = H(i) + hk
   -1
108             - provide nonlinearity and translation
   invariance
109             - global average pooling before FC to
   reduce computation load
110         - Skip connection
111             - provide unimpeded gradient flow
112             - provide multiple level of
   abstractions and let the network itself to decide
   which level to use; the network
113                 becomes a "bag" of different models,
   similar to ensemble
114             - the above point can be seen from an
   optimization perspective, in which says that
   deeper networks have more
115                 complicated loss surface and require
   much more time and more sophisticated optimization
   techniques to converge.
116                 The skip connections ease this
   difficulty by allowing the network to converge at
   "less-representative" minima
117
```

```
118  7. CNN - training
119      - preprocessing
120          - zero-centered
121          - normalization
122      - weight initialization
123          - small random number (gaussian with zero
     mean and 1e-2 standard deviation)
124              - for deeper networks, activation
     outputs become zero
125              - weight updating becomes super slow,
     sometimes completely stop, G[l-1] = X^T * G[l-1]
126          - Xavier initialization
127              - small random / sqrt(fan_in)
128      - batch normalization
129          - covariate shift: the change of the
     distribution of input data
130          - almost eliminate gradient vanishing,
     alleviate internal covariate shift, regularization
     , network converge faster, can use larger learning
      rate
131          - at test time, should use empirical
     parameters obtained at training stage
132          - for fc layer, we do per-dimension batch
     norm; for conv, we do per-channel batch norm
133      - optimization
134          - problems with sgd:
135              - stuck in local minima
136              - slow at saddle point
137          - momentum: v[t] = alpha * v[t-1] + G; w[t
     +1] = w[t] - lr * v[t]
138              - jump over local minima
139              - speed up at saddle point
140          - second-order optimization
141              - learning rate determined by hessian
     matrix, point to minima so converge faster
142              - computationally expensive
143      - model ensembles
144          - train multiple independent model
145          - at test time, take the average of their
     results
```

```
146
```