

ECE/CS 559 Lecture 13 T 10/8

Last time: Regularization

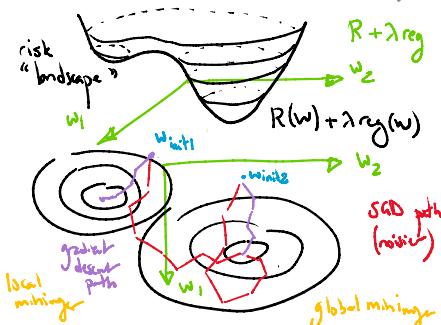
- Idea: Avoid overexplaining by limiting to reasonable explanations: prior knowledge / constraints on weights/biases
- $\min_w R(w) \text{ s.t. } \|w\|_2^2 \leq C \iff \min_w R(w) + \lambda \|w\|_2^2 \quad L_2$
- During gradient descent simply add: $\nabla_w (R(w) + \lambda \|w\|_2^2) = \nabla_w R + 2\lambda w$ (other choice)
- Dropout: Randomly prune weights during training (with prob. p)
Include all weights (scaled by p) at test time.

• Reasons for minibatch SGD:

- Computation: It's faster to get an update with small batches
- Statistics: Heterogeneous data \Rightarrow small batches give a glimpse of the whole
- Optimization: The variance can be helpful to get out of local minima.
(choose batch size to trade off the pert with convergence speed.)

• General form of SGD:

for epoch in # of epochs:
 $t=0; \Sigma \nabla = 0$
 For (key) in Data:
 Forward(x); backward(y)
 $\Sigma \nabla = \nabla R + \lambda \nabla g_y; t=t+1$
 $\text{if } (\nabla \cdot \nabla) \cdot \nabla = 0:$
 $w += -\eta \Sigma \nabla$
 $\Sigma \nabla = 0$



② Heuristics The "craft" of training neural networks.

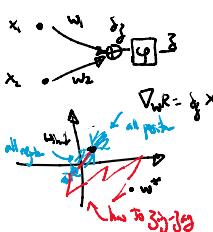
① Choosing the activation function:

- For the longest time: sigmoid s , hyperbolic tangent $tanh$
- Issue: $(s'(x))^2 \rightarrow$ saturation, slow down learning



② Preprocessing inputs:

- Ideally, input coordinates are uncorrelated/centered
 - Otherwise, may slow down learning
 - Example: if x_i 's are always + or - all together
- Preprocess to make them have 0 mean, unit variance, 0 covariance: "whitened"



① Optimization

② Stochastic Gradient Descent:

Empirical risk $R(w) = \frac{1}{n} \sum_{(y, x) \in \text{Data}} l(y, f(x; w))$
 "True" risk $R'(w) = \mathbb{E} [l(Y, f(X; w))]$

Note: $R(w)$ is an unbiased estimate of $R'(w)$: $\mathbb{E}[R(w)] = R'(w)$

This means the gradients are unbiased: $\mathbb{E}[\nabla_w R(w)] = \nabla_w R'(w)$

- But... this remains true for even a single point $\mathbb{E}[\nabla_w l(y, f(x; w))] = \nabla_w R(w)$
- or a mini batch of points $\mathbb{E}[\nabla_w \sum_{(y, x) \in B} l(y, f(x; w))] = \nabla_w R(w)$
- What changes? The variance. # points \downarrow variance \uparrow noisiness \uparrow
Then why don't we always use all the points?

③ Choice of η , the "learning rate"

- Because SGD introduces noise, η should be adjusted
 - Otherwise, performance can worsen/improve/worsen (oscillate)
- Typically: Reduce η by a factor (e.g. 90%) if sustained worsening.
- Advanced: Adapt to the landscape (e.g. flatness, steepness)
- Vary η by layer: Later layers $\nabla_w \uparrow$, make $\eta \downarrow$
by neuron: # inputs \uparrow , make $\eta \downarrow \propto 1/\sqrt{\# \text{inputs}}$ (heuristic)

④ Controlling instability:

- Gradients are repeatedly multiplied during backpropagation.
 - As a result, they could vanish ($\times 0.9 \times 0.9 \times \dots$) or explode ($\times 1 \times 1 \times \dots$)
 - Solution: normalize gradients (across a batch or a layer), clip gradients

mean (vector) $\mu = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} x$ covariance (matrix) $C = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} (x - \mu)(x - \mu)^T$
 C is positive semidefinite: can be diagonalized $C = UDU^T$, $UU^T = U^T U = I$
 Let $x_{\text{new}} = D^{-\frac{1}{2}} U^T (x - \mu)$ \leftarrow has 0 mean, covariance = I

⑤ Initializing the weights:

- Intuition: maintain similar statistics from layer to layer.
- Ignoring activation, assuming whitened inputs, covariance I , next layer's inputs is:
 $C_i = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} (Wx)(Wx)^T = WW^T = n \begin{bmatrix} & \\ & \ddots \\ & & n \end{bmatrix}$
- Ideally, we want to make $\mathbb{E}[WW^T] = I$
- If w have 0 mean + uncorrelated, then off-diagonals are 0 ✓
- On the diagonal, we'd be adding n , $\mathbb{E}[w^2]$, so we get $n \text{ var}(w)$.
- To make it identity, we just have to choose $\text{var}(w) = 1$
- E.g., choose w to be i.i.d. $N(0, \frac{1}{\sqrt{n}})$