

Last time: Chain Rule, logistic regression

$$p(y=1|x;w) = \frac{1}{1 + \exp\{-f(x,w)\}}$$

$$p(y=0|x;w) = \frac{1}{1 + \exp\{f(x,w)\}}$$

$$\mathcal{L}(w) = \sum_n y_n \ln(1 + \exp\{-f_n\}) + (1 - y_n) \ln(1 + \exp\{f_n\})$$

notice: could plug in other expressions for $p(y=1|x)$, $p(y=0|x)$
wrote loss in terms of f_n ,

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial f_n} \frac{\partial f_n}{\partial w} \quad \text{e.g. } f_n = w^T x_n + w_0$$

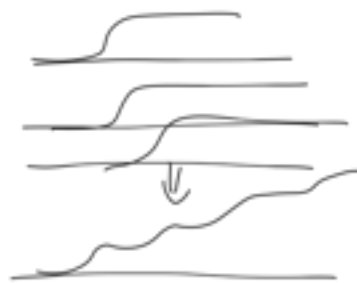
Applied This idea to shallow network (basis regression)



$$\phi_j = \sigma(w_j^{*T} x + w_j^+)$$

nonlinearity \Rightarrow sigmoid/logit; ReLU/hinge; tanh

deep networks:



How are we going to fit these models? $(W^k, w^k, w) \Rightarrow$ lots of params!

problem: very non-convex! (can be a real issue \rightarrow in practice:

lots of things to make the problem easier \Rightarrow random restarts,
may use various forms of momentum/stochasticity, overparameterization)

How to opt. NNs with SGD, key thing we need is the gradient of the loss w.r.t. params.

$$\text{regression loss: } \mathcal{L}(w, W) = \frac{1}{2} \sum_n (y_n - f_n)^2$$

$$\frac{\partial \mathcal{L}}{\partial f_n} = \boxed{-(y_n - f_n)}$$

Aside: vector-valued chain rule

$$\text{scalar chain rule: } y = f(u) = f(g(h(x)))$$

$$u = g(v)$$

$$v = h(x)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial u} \frac{\partial g}{\partial v} \frac{\partial h}{\partial x} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial v} \frac{\partial v}{\partial x}$$

now, suppose x was a vector. $\Rightarrow \nabla_x f = \frac{\partial f}{\partial u} \frac{\partial u}{\partial v} \underbrace{\nabla_x v}_{\text{vector}}$

if v was also a vector:

$$\nabla_x f = \underbrace{\frac{\partial f}{\partial u}}_{\text{scalar}} \underbrace{\nabla_u [u]}_{\text{vector of size } J} \underbrace{J_x [v]}_{\text{matrix } J \times D}$$

$$x \in \mathbb{R}^D \quad v \in \mathbb{R}^J$$

if u were also a vector:

$$\nabla_x f = \nabla_u f J_u [u] J_x [v]$$

Now we're ready to start optimizing. Let's start w/w:

$$x \xrightarrow{W^1} \phi \xrightarrow{W} f$$

$$\nabla_w \mathcal{L}_n = \underbrace{\frac{\partial \mathcal{L}_n}{\partial f_n}}_{-(y_n - f_n)} \underbrace{\nabla_{f_n} f_n}_{\nabla_{f_n} f_n = \phi_n} = \underbrace{-(y_n - f_n)}_{\text{scalar}} \underbrace{(\phi_n)}_{\text{vector of size } J}$$

Notice! We need f_n to compute $\nabla_{f_n} \mathcal{L}_n$, and f_n depends on w . calc feed-forward (values)

- 1) for our current w , we compute f_n
- 2) which allows us to compute $\nabla_{f_n} \mathcal{L}_n$ w.r.t. our current w . backward pass (calc gradients)

$$\nabla_{W^1} \mathcal{L}_n = \underbrace{\frac{\partial \mathcal{L}_n}{\partial f_n}}_{-(y_n - f_n)} \underbrace{\nabla_{f_n} f_n}_{\nabla_{f_n} f_n = \phi_n} \underbrace{J_{W^1} [\phi_n]}_{J \times J^T}$$

lets flatten W^1 into \tilde{W}_{dj}^1

$$\phi_{nj} = \sigma(\sum_d \tilde{W}_{dj}^1 x_d + w_j^1)$$

$$\frac{\partial \phi_{nj}}{\partial \tilde{W}_{dj}^1} = 0 \text{ if } j \neq j'$$

$$= \sigma(\sum_d \tilde{W}_{dj}^1 x_d + w_j^1) \cdot (1 - \sigma(\sum_d \tilde{W}_{dj}^1 x_d + w_j^1)) \cdot x_d$$

we get a $J \times JD$ matrix

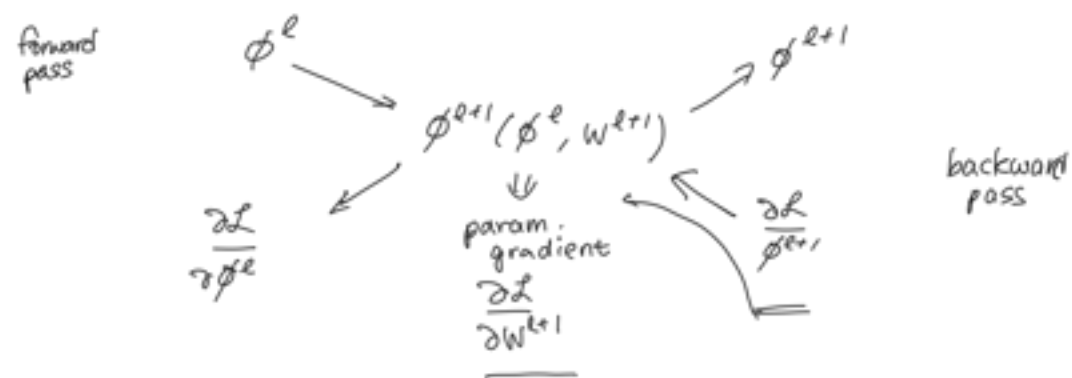


General approach: compute $\mathcal{L}(\phi^L(\phi^{L-1}(\dots \phi^1(x))))$ for L layers
store all intermediate values ϕ^e
 (forward pass)

$$\text{Backward pass: compute } \frac{\partial \mathcal{L}}{\partial \phi^L} = \nabla_{\phi^L} \mathcal{L} \underbrace{J_{\phi^{L-1}} [\phi^L]}_{\text{Jacobian}} \underbrace{J_{\phi^{L-2}} [\phi^{L-1}]}_{\text{Jacobian}} \dots \underbrace{J_{\phi^1} [\phi^2]}_{\text{Jacobian}}$$

$$\text{if we want } \frac{\partial \mathcal{L}}{\partial \phi^{L-1}} = \frac{\partial \mathcal{L}}{\partial \phi^L} J_{\phi^{L-1}} [\phi^L]$$

$$\text{Compute parameter gradients } \frac{\partial \mathcal{L}}{\partial W^L} = \nabla_{\phi^L} \mathcal{L} J_{\phi^{L-1}} [\phi^L] \dots J_{W^L} [\phi^L]$$



computation graph

What's special about NNs now?

- ⇒ we have more data now
- ⇒ we have GPUs / software tools for autodiff
- ⇒ that said, lots of practicalities

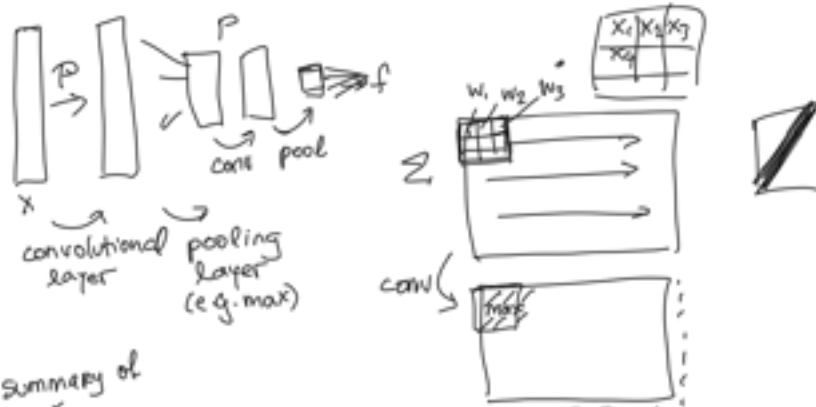
- ~~overfitting~~ overfitting:
 - L1/L2
 - perturbations on inputs
 - dropout
- choices of activation (sigmoid, tanh, etc.) to avoid saturation
- SGD: ways of variance reduction / subgradients / ADAM

Different Architectures:

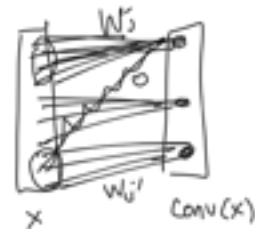
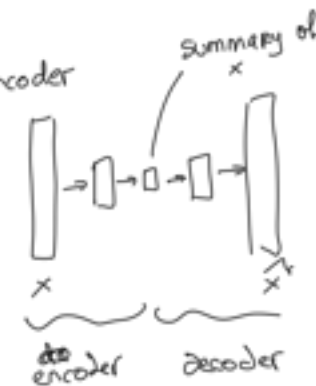
Feedforward Network,
Multi-layer Perceptron



Convolutional NN (parameter tying)



Autoencoder

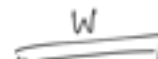


RNN: Recurrent NN:



(variant: long short-term memory - LSTM)
gated recurrent units GRU

Concept exercise:



Suppose we have an image that is 

- 1) how many params would be reqd in FF NN? Assume L layers, size M nodes, single output. $(DN)(M) + \underbrace{M^2(L-1)}_{\text{ignore for now}} + M$
- 2) conv net w/ patch size P , L layers (fully connected single output)

$$\underbrace{(LP^2)}_{\text{ignore for now}} + \underbrace{(W-LP)(D-LP)}_{\text{ignore for now}} / \text{some pooling factor}$$

- 3) RNN that travels from $d=0$ to $d=D$ w/ single output, H = dim hidden layer?
- $1, 2, \dots, \dots$