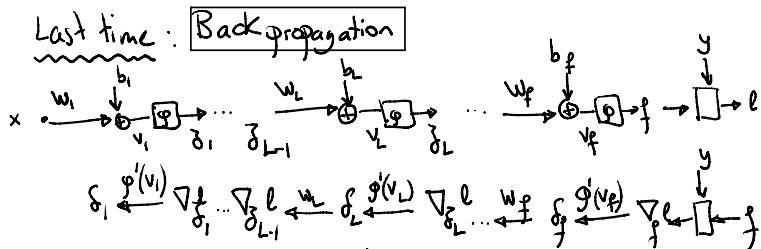


L11

ECE/CS 559 Lecture 11

T 10/1

Last time: Backpropagation

$$\text{Forward: } l(\delta_L) = l(g(W_L \delta_{L-1} + b_L)) \quad \text{element wise}$$

$$\text{Backward: } \nabla_{\delta_{L-1}} l = W_L^T (g'(v_L) \cdot \nabla_{\delta_L} l) = W_L^T \delta_L$$

$$\text{Gradients: } \nabla_{W_L} l = (g'(v_L) \cdot \nabla_{\delta_L} l) \delta_L^T = \delta_L \delta_L^T, \quad \nabla_{b_L} l = \delta_L$$

- When learning neural networks, issue is that we train on the observed contexts x . However, we mostly care about performance on new unobserved contexts!
- Performing well in unobserved contexts = generalization
overfitting prevents generalization (unless it's mild)
- How do we quantify this?

Training Data = $\{(x, y), \dots\}$ Testing Data = $\{(x', y'), \dots\}$
 ↳ use this to build NN ↳ use this to assess generalization

We must be careful!

We can't keep on using testing data to refine the NN, otherwise we may overfit it too and not generalize to new data.

(2) Cross-Validation

- Since training risk is not a good enough proxy to testing risk, we could hold out some of the training data for validation:

Training Data = $\{(x, y), \dots\}$ Validation Data = $\{(x', y'), \dots\}$, Testing Data = $\{(x'', y'')\}$
 ↴ train on this (open) ↴ check if okay on this (green) ↴ test on this (blue)

- What is validation good for?
 - Choose what epoch to stop at.
 - Choose between possible architectures, learning rates, etc.
 - These choices are "coarser" than choosing the weights & biases

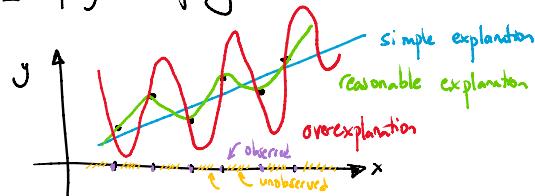
k-Fold Cross-validation: 

Change validation block, repeat. Average.

Note: Any change that helps prefer some w over others is regularization. This introduces what we call inductive bias, which (if right) improves generalization.

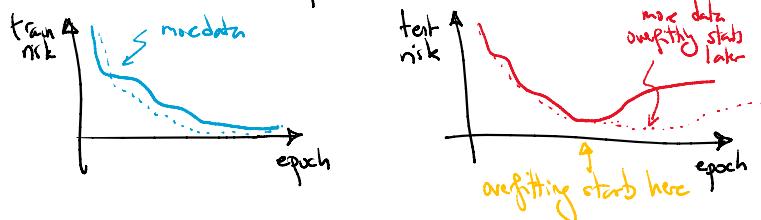
(1) Overfitting

- Tendency to overexplain the training data.
- Overexplain = when many explanations describe experience equally well, choose one that adds many complicated details and twists (typically what conspiracy theories do)
- Example: polynomial fitting



- Typical property leading to this: high sensitivity to changes in data

(2) Typical behavior of train vs. test risk:



(3) Observations:

- Stopping early could help alleviate overfitting
Why? The NN will not fit the noise.
How? By staying closer to its initialization (staying "simpler")
- More training data reduces overfitting, delays its onset.

(3) Regularization

- Cross-validation helps see if we overfit. Regularization helps prevent it.
- The word comes from inverse problems in physics:

- To decide between many possible explanations, we need to favor some.
- We tend to favor "simple" ones (Occam's Razor) (e.g., low energy)
- E.g., for polynomials and NNs, favor small weights:

$$\min_w R(w) \text{ for } \|w\|^2 \leq C \iff \min_w R(w) + \lambda \|w\|_2^2$$

regularized/penalized risk
- What does it mean? Prefer smoother functions: $\|w\|_1 = |w_1| + \dots + |w_n|$ (other choice)
- Why does it work? Can't overfit noise

- Regularized risk is a better proxy for test-risk
- How does it affect backpropagation? Simple!

$$\nabla_w (R(w) + \lambda \|w\|_2^2) = \nabla_w R + 2\lambda w \underbrace{\text{sign}(w)}_{\|w\|_1}$$

L9&10

ECE/CS 559 Lecture 9 & 10 T,Th 9/24, 9/26

Last time: Gradient Descent, Widrow-Hoff LMS Algorithm

$$\text{minimize } R(w) = \sum_{(x,y) \in \text{Data}} l(y, f(x; w))$$

$$w' \leftarrow w - \eta \nabla R(w)$$

$$\nabla R_w = \left[\frac{\partial R}{\partial w_1}, \dots, \frac{\partial R}{\partial w_k} \right]$$

$$\cdot l(y, f(x; w)) = \|y - f(x; w)\|^2$$

$$f(x; w) = w^T x$$

$$w' \leftarrow w - \eta \nabla R(w) = w + 2\eta \sum_{(x,y)} (y - w^T x) x^T$$

- We will apply the chain rule multiple times:

$$\frac{\partial R}{\partial w_k} = \sum_{(x,y)} \left[\frac{\partial l}{\partial w_k} \right]^{(*)}$$

need this

$$(?) \frac{\partial l}{\partial w_k} = \frac{\partial}{\partial w_k} l(t_k, \dots) \Big|_{t_k = g(w_k z + \dots)} \quad \begin{matrix} v_{t_k} \\ \vdots \\ \end{matrix}$$

have this from forward pass

$$(\text{two chain rules}) = \left[\frac{\partial e}{\partial t_k} \Big|_{t_k = g(v_{t_k})} \right]^{(*)} \cdot g'(v_{t_k}) \cdot z$$

need this at every neuron output

$$(*) \frac{\partial l}{\partial z} = \frac{\partial}{\partial z} l(t_1, \dots, t_k, \dots, t_m) \Big|_{t_k = g(w_k z + \dots)} \quad \begin{matrix} v_{t_k} \\ \vdots \\ \end{matrix}$$

(two chain rules)

$$= \sum_{k=1}^m \frac{\partial l}{\partial t_k} \Big|_{t_k = g(v_{t_k})} \cdot g'(v_{t_k}) \cdot w_k$$

- Of course, the layer of z has other (say n) neurons. Let's index the equations for these (blue j)

$$\frac{\partial l}{\partial w_{kj}} = \delta_{t_k} \cdot z_j$$

$$\frac{\partial l}{\partial z_j} = \sum_{k=1}^m \delta_{t_k} w_{kj} \quad j=1, \dots, n$$

- If we think of δ_t as a m -dim vector of w as a $m \times n$ -dim matrix
- We get the linear algebraic backward equation:

$$\nabla_l = \delta_{t_k} z^T$$

outer product

$$\nabla_l = W^T \delta_t$$

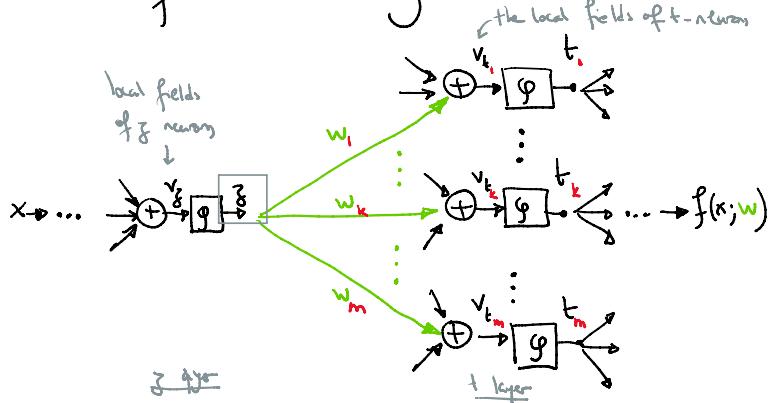
matrix multiplication

$$\delta_z = \nabla_l \circ g'(v_z)$$

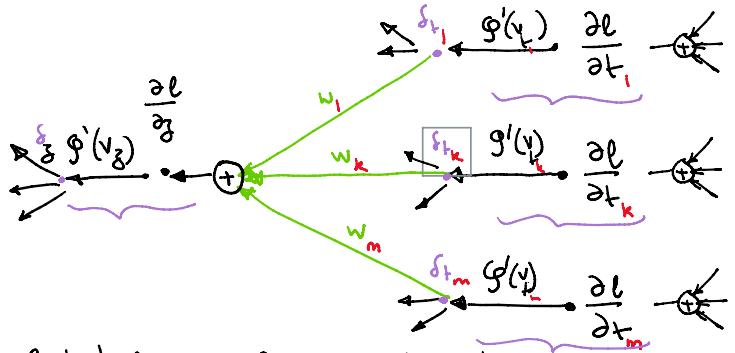
elementwise product

Back propagation

- So far we've been calculating gradients with 1 neuron.
- What if we connect many neurons?



- To calculate (?) and (*), we perform a forward pass to get all the outputs (v_{t_1}, v_{t_2}, \dots), then we perform a backward pass:

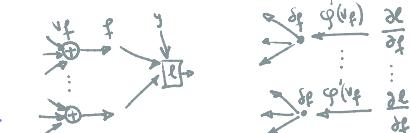


Gradient of w = forward signal \times backward signal δ_t

- To complete this, we need 2 things: how to start, how to handle bias

- Start at the last layer ∇_l depends on the choice of l

e.g. if $l = \|f - y\|^2$ then
 $\nabla_l = 2(f - y)$, $\delta_f = \nabla_l \circ g'(v_f)$



- Biases are like dead-ends on the backward path:

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial b} \Big|_{t_k = g(v_k + b)} \quad \begin{matrix} v_k \\ \vdots \\ \end{matrix}$$

$$= \frac{\partial l}{\partial t_k} \Big|_{t_k = g(v_k)} \cdot g'(v_k) \cdot 1$$

$\nabla_l = \delta_t$

$$\nabla_l = \delta_t \circ g'(v_t) \quad \begin{matrix} v_t \\ \vdots \\ \end{matrix}$$

$$\nabla_b l = \frac{[\nabla_l \circ g'(v_t)] \delta_t}{\delta_t}$$

$$\nabla_z l = W^T [\nabla_l \circ g'(v_z)]$$

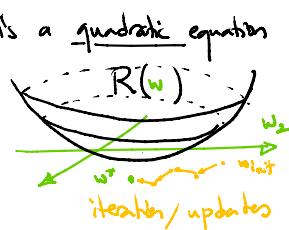
Note: Can we find eqs to get these: $l(t) = l(\underbrace{g(w_j + b)}_t)$ $\Rightarrow \nabla_w l = \frac{[\nabla_l \circ g'(v_t)] \delta_t}{\delta_t}$

L8

ECE/CS 559 Lecture 8 Th 9/19

Last time: Linear Regression with Squared Loss

- Data: $\{(x_i, y_i)\}_{i=1}^n$ $x \in \mathbb{R}^m$, $y \in \mathbb{R}$
- Predictor: $f(x; w) = w^T x$ $w \in \mathbb{R}^m$
- Goal: minimize $R(w) = \sum_{(x_i, y_i)} \|y_i - f(x_i; w)\|^2$
- $w = Y [X^T (X X^T)^{-1}]$ Pseudo inverse



1-D Example: $R(w) = w^4 - 4w^3 + 3w^2$ $\eta = \frac{1}{8}$

$$\begin{aligned} w(0) &= 1 & w(1) &\leftarrow w(0) - \nabla R(w(0)) \cdot \eta \\ && 1 &- 4 \cdot 1^3 \cdot \frac{1}{8} = \frac{1}{2} \\ w(2) &\leftarrow w(1) - \nabla R(w(1)) \cdot \eta \\ &\frac{1}{2} &- 4 \left(\frac{1}{2}\right)^3 \cdot \frac{1}{8} = \frac{3}{16} \end{aligned}$$

What's happening?

- $w(t) \leftarrow w(t-1) - 4 w(t-1)^3 \eta = [1 - 4\eta w(t-1)^2] w(t-1)$
- If η is small enough, always diminishes. (< 1) Since bounded from below (by 0) will converge.
- But will it converge to 0? Yes, think about why.
- What if η is too big? Try $\eta = \frac{1}{2}$ with $w(0) = 1$.

(2) Delta Rule

Squared loss for single neuron with differentiable $g(\cdot)$:

$$\begin{aligned} R(w) &= \sum_{(x_i, y_i)} (y_i - g(w^T x_i))^2 \\ \nabla R(w) &= -2 \sum_{(x_i, y_i)} (y_i - g(w^T x_i)) g'(w^T x_i) x_i \\ w &\leftarrow w + 2\eta \sum_{(x_i, y_i)} (y_i - g(w^T x_i)) g'(w^T x_i) x_i \end{aligned}$$

1 data point at a time:

$$w \leftarrow w + \eta (y - g(w^T x)) g'(w^T x) x$$

(1) Gradient Descent

$$\text{Goal: minimize } R(w) = \sum_{(x_i, y_i) \in \text{Data}} \|y_i - f(x_i; w)\|^2$$

Idea: Set the derivative to zero!

Gradient: $\nabla R_w = [\frac{\partial R}{\partial w_1}, \dots, \frac{\partial R}{\partial w_k}]$ reduce step-by-step.

Each step: Change w by Δw : $w' \leftarrow w + \Delta w$.

$$R(w') = R(w + \Delta w) = R(w) + \nabla R(w)^\top \Delta w + O(\|\Delta w\|^2)$$

How should we choose Δw to get the best reduction?

$$\text{To minimize } \nabla R^\top \Delta w, \text{ let } \Delta w \propto -\nabla R: w' \leftarrow w - \eta \nabla R$$

(2) Widrow-Hoff LMS Algorithm

Let's apply gradient descent to linear regression. Recall:

$$\nabla R(w) = -2 \sum_{(x_i, y_i) \in \text{Data}} (y_i - w^T x_i) x_i^\top \quad \text{same shape as } W$$

$$w \leftarrow w - \eta \nabla R(w) = w + 2\eta \sum_{(x_i, y_i)} (y_i - w^T x_i) x_i^\top$$

$$\text{When } y \in \mathbb{R}^1 \text{ then } w = w^\top \text{ and } y - w^\top x \in \mathbb{R}$$

$$\text{The update becomes: } w \leftarrow w + \eta \sum_{(x_i, y_i)} (y_i - w^\top x_i) x_i$$

We could do this 1 data point at a time:

$$\text{For each } (x_i, y_i) \in \text{Data}: w \leftarrow w + \eta (y_i - w^\top x_i) x_i$$

Similar to the perceptron learning algorithm!

Example $g(\cdot)$:

$$\text{Sigmoid: } g(v) = \frac{1}{1 + e^{-av}} \quad \text{since } 1 - g(v) = \frac{e^{-av}}{1 + e^{-av}}$$

$$g'(v) = \frac{0 - (-ae^{-av})}{(1 + e^{-av})^2} = \frac{ae^{-av}}{(1 + e^{-av})^2} = a g(v) (1 - g(v))$$

$$\text{ReLU: } g(v) = \begin{cases} v & v \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$g'(v) = \begin{cases} 1 & v \geq 0 \\ 0 & \text{otherwise} \end{cases} = \text{step}(v)$$

L7

ECE/CS 559 Lecture 7

9/17

Last time : Perception Learning Example Supervised Learning

- Data: $\{(x_i, y_i), \dots\}$ $x \in \mathcal{X}$ (feature space), $y \in \mathcal{Y}$ (label space)
- Task: Predict/initiate y using only x .
- How: Use a predictor neural network $y_{\mathbf{w}}(x) = f(x; \mathbf{w})$
- Loss: how far y and $f(x; \mathbf{w})$ are: $\ell(y, f) = \frac{1}{2} \{y - f\}^2$
- Risk: (empirical) $R(\mathbf{w}) = \frac{1}{|\text{Data}|} \sum_{(x_i, y_i) \in \text{Data}} \ell(y_i, f(x_i; \mathbf{w}))$
- Goal: Minimize (empirical) risk,
make few mistakes/errors / stay close to y .

• Types of algorithms :

- See (x_i, y_i) , update \mathbf{w}
- See (x_i, y_i) , update \mathbf{w}
- See (x_i, y_i) , update \mathbf{w}

- Get Data = $\{(x_i, y_i), \dots\}$
 Epoch 1: use Data, update \mathbf{w}
 Epoch 2: use Data, update \mathbf{w}



① Linear Regression with Squared Loss

- Data: $\{(x_i, y_i), \dots\} \quad x \in \mathbb{R}^n \quad y \in \mathbb{R}^m$
- Predictor: $f(x; \mathbf{w}) = \mathbf{W}x \quad \mathbf{W} \in \mathbb{R}^{m \times n}$
(this is a neuron with $g(v) = v$, identity activation)
- Goal: minimize $R(\mathbf{w}) = \sum_{(x_i, y_i) \in \text{Data}} \|y_i - \mathbf{W}x_i\|^2$

Analytic Solution: Optimality conditions, unconstrained $\Rightarrow \nabla_{\mathbf{w}} R = 0$

$$\begin{aligned} \cdot \text{Single data point: } l &= \sum_{k=1}^m (y_k - \sum_{j=1}^n w_{kj} x_j)^2 \quad \text{gradient} \\ \frac{\partial l}{\partial w_{kj}} &= -2x_j (y_k - \sum_{j=1}^n w_{kj} x_j) = -2y_k x_j + \sum_{j=1}^n w_{kj} (x_j x_j) \\ \downarrow \begin{bmatrix} 1 & j & n \end{bmatrix} \rightarrow \nabla_{\mathbf{w}} l &= -2y x^T + \mathbf{W} x x^T \\ \cdot \text{All data points: } R(\mathbf{w}) &= \sum_{(x_i, y_i) \in \text{Data}} \ell(y_i, \mathbf{W}x_i) \Rightarrow \nabla_{\mathbf{w}} R = \sum_{(x_i, y_i) \in \text{Data}} \nabla_{\mathbf{w}} \ell \\ \nabla_{\mathbf{w}} R &= \sum_{(x_i, y_i) \in \text{Data}} -2y x^T + 2\mathbf{W} \sum_{x \in \text{Data}} x x^T \\ &= -2 \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^T}_{\text{Data}} + 2\mathbf{W} \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^T}_{\text{Data}} \\ &= -2 \mathbf{Y} \mathbf{X}^T + 2\mathbf{W} \mathbf{X} \mathbf{X}^T \end{aligned}$$

$$\bullet \nabla_{\mathbf{W}} R = 0 \Rightarrow -2 \mathbf{Y} \mathbf{X}^T + 2 \mathbf{W} \mathbf{X} \mathbf{X}^T = 0$$

$$\Rightarrow \mathbf{W} = \mathbf{Y} \boxed{\mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1}}$$

Moore-Penrose
Pseudo inverse

$$\text{Compare to: } \mathbf{W} \mathbf{X} \approx \mathbf{Y} \quad \text{not invertible if } |\text{Data}| > n$$

$$\mathbf{W} \begin{bmatrix} \mathbf{X} \\ \vdots \end{bmatrix} \approx \mathbf{Y} \quad \text{but it's "like" invertible if above.}$$

• Where are the updates? One single batch update!

Intuition: it's a quadratic equation



• Other shapes → solve by iteration/updates

② Gradient Descent

Goal: minimize $R(\mathbf{w}) = \sum_{(x_i, y_i) \in \text{Data}} \|y_i - f(x_i; \mathbf{w})\|^2$

Idea: Set the derivative to zero!

Gradient: $\nabla_{\mathbf{w}} R = \left[\frac{\partial R}{\partial w_1}, \dots, \frac{\partial R}{\partial w_k} \right]$ reduce step-by-step.

Each step: Change w by Δw : $w' \leftarrow w + \Delta w$.

$$R(w) = R(w + \Delta w) = R(w) + \nabla_{\mathbf{w}} R^T \Delta w + O(\|\Delta w\|^2)$$

Taylor expansion ignore

How should we choose Δw to get the best reduction?

- To minimize $\nabla_{\mathbf{w}} R^T \Delta w$, let $\Delta w \propto -\nabla_{\mathbf{w}} R$: $w' \leftarrow w - \eta \nabla_{\mathbf{w}} R$
- ↑ proportional

$$1\text{-D Example: } R(w) = w^4 \quad \frac{dR}{dw} = 4w^3 \quad \eta = \frac{1}{8}$$

$$\begin{aligned} w(0) &= 1 & w(1) &\leftarrow w(0) - \nabla R(w(0)) \cdot \eta \\ && 1 &- 4 \cdot \frac{1}{8} = \frac{1}{2} \\ w(2) &\leftarrow w(1) - \nabla R(w(1)) \cdot \eta \\ &\frac{1}{2} &- 4 \left(\frac{1}{2} \right)^3 \cdot \frac{1}{8} = \frac{3}{16} \end{aligned}$$

• What's happening?

$$w(t) \leftarrow w(t-1) - 4 w(t-1)^3 \eta = [1 - 4\eta w(t-1)^2] w(t-1)$$

- If η is small enough, always diminishes. (< 1) Since bounded from below (by 0) will converge.

- But will it converge to 0? Yes, think about why.

- What if η is too big? Try $\eta = \frac{1}{2}$ with $w(0) = 1$.