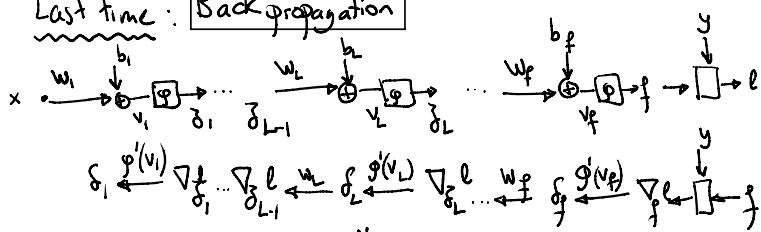


## ECE/CS 559 Lecture 11

T 10/1

Last time: Backpropagation

$$\text{Forward: } l(\delta_L) = l(g(\underbrace{w_L \delta_{L-1} + b_L}_{v_L})) \text{ element wise}$$

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

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

- 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 =  $\{(x_i, y_i), \dots\}$  Testing Data =  $\{(x'_i, y'_i), \dots\}$   
 ↳ 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 =  $\{(x_i, y_i), \dots\}$  Validation Data =  $\{(x'_i, y'_i), \dots\}$ , Testing Data =  $\{(x''_i, y''_i)\}$   
 train on this (green) check if okay on this (yellow) 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

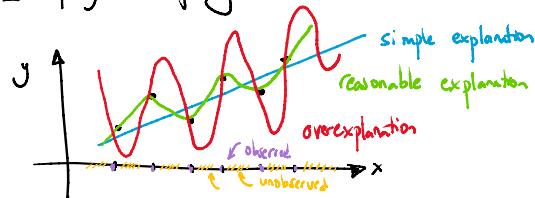
k-Fold Cross-validation:

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.

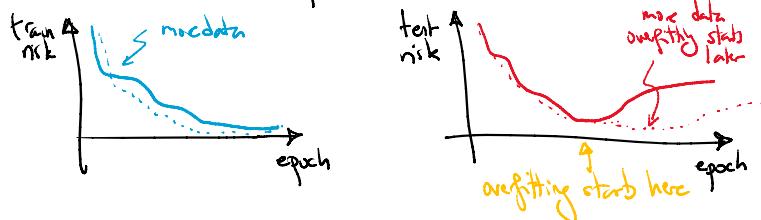
## (1) Overfitting

- 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!

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

$|w|$