

# CSC321 Lecture 9: Generalization

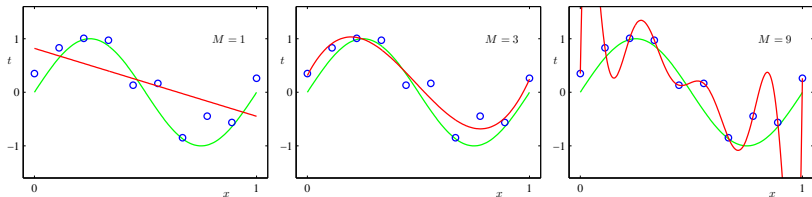
Roger Grosse

# Overview

- We've focused so far on how to *optimize* neural nets — how to get them to make good predictions on the training set.
- How do we make sure they generalize to data they haven't seen before?
- Even though the topic is well studied, it's still poorly understood.

# Generalization

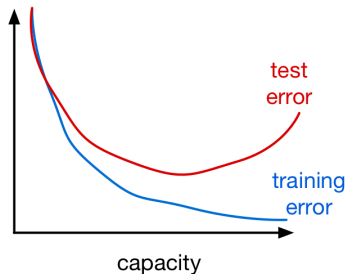
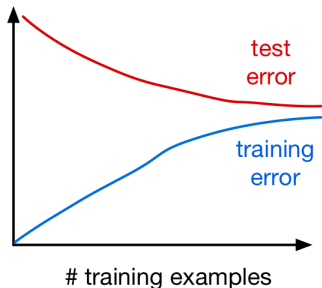
## Recall: overfitting and underfitting



We'd like to minimize the generalization error, i.e. error on novel examples.

# Generalization

- Training and test error as a function of data and model capacity.
- Note: capacity isn't a formal term, or a quantity we can measure.



# Bias/Variance Decomposition

- There's an interesting decomposition of generalization error in the particular case of squared error loss.
- It's often convenient to suppose our training and test data are sampled from a **data generating distribution**  $p_{\mathcal{D}}(\mathbf{x}, t)$ .
- Suppose we're given an input  $\mathbf{x}$ . We'd like to minimize the expected loss:

$$\mathbb{E}_{p_{\mathcal{D}}}[(y - t)^2 | \mathbf{x}]$$

- The best possible prediction we can make is the conditional expectation

$$y_{\star} = \mathbb{E}_{p_{\mathcal{D}}}[t | \mathbf{x}].$$

- Proof:

$$\begin{aligned}\mathbb{E}[(y - t)^2 | \mathbf{x}] &= \mathbb{E}[y^2 - 2yt + t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t | \mathbf{x}]^2 + \text{Var}[t | \mathbf{x}] \\ &= (y - y_{\star})^2 + \text{Var}[t | \mathbf{x}]\end{aligned}$$

- The term  $\text{Var}[t | \mathbf{x}]$ , called the **Bayes error**, is the best risk we can hope to achieve.

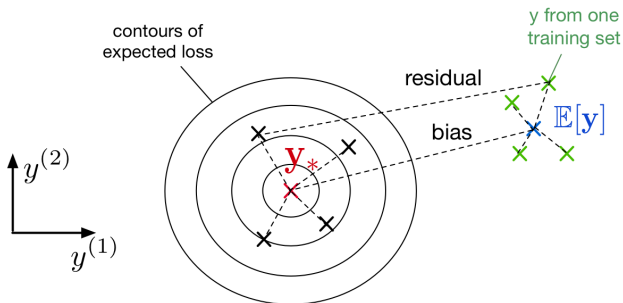
# Bias/Variance Decomposition

- Now suppose we sample a training set from the data generating distribution, train a model on it, and use that model to make a prediction  $y$  on test example  $\mathbf{x}$ .
- Here,  $y$  is a random variable, and we get a new value each time we sample a new training set.
- We'd like to minimize the **risk**, or expected loss  $\mathbb{E}[(y - t)^2]$ . We can decompose this into **bias**, **variance**, and Bayes error. (We suppress the conditioning on  $\mathbf{x}$  for clarity.)

$$\begin{aligned}\mathbb{E}[(y - t)^2] &= \mathbb{E}[(y - y_*)^2] + \text{Var}(t) \\ &= \mathbb{E}[y_*^2 - 2y_*y + y^2] + \text{Var}(t) \\ &= y_*^2 - 2y_*\mathbb{E}[y] + \mathbb{E}[y^2] + \text{Var}(t) \\ &= y_*^2 - 2y_*\mathbb{E}[y] + \mathbb{E}[y]^2 + \text{Var}(y) + \text{Var}(t) \\ &= \underbrace{(y_* - \mathbb{E}[y])^2}_{\text{bias}} + \underbrace{\text{Var}(y)}_{\text{variance}} + \underbrace{\text{Var}(t)}_{\text{Bayes error}}\end{aligned}$$

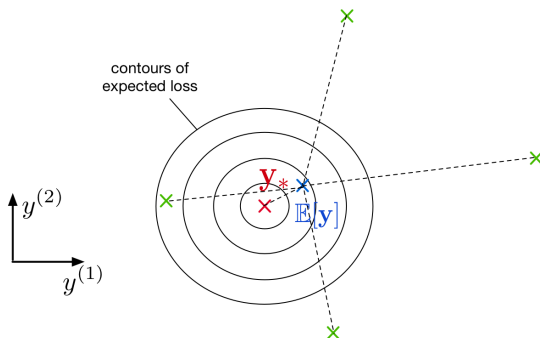
# Bias/Variance Decomposition

- We can visualize this decomposition in **prediction space**, where the axes correspond to predictions on the test examples.
- If we have an overly simple model, it might have
  - high bias (because it's too simplistic to capture the structure in the data)
  - low variance (because there's enough data to estimate the parameters, so you get the same results for any sample of the training data)



# Bias/Variance Decomposition

- If you have an overly complex model, it might have
  - low bias (since it learns all the relevant structure)
  - high variance (it fits the idiosyncrasies of the data you happened to sample)



- The bias/variance decomposition holds only for squared error, but it provides a useful intuition even for other loss functions.



# Our Bag of Tricks

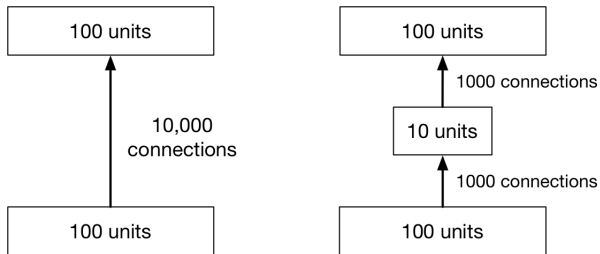
- How can we train a model that's complex enough to model the structure in the data, but prevent it from overfitting? I.e., how to achieve low bias and low variance?
- Our bag of tricks
  - data augmentation
  - reduce the capacity (e.g. number of paramters)
  - weight decay
  - early stopping
  - ensembles (combine predictions of different models)
  - stochastic regularization (e.g. dropout, batch normalization)
- The best-performing models on most benchmarks use some or all of these tricks.

# Data Augmentation

- The best way to improve generalization is to collect more data!
- Suppose we already have all the data we're willing to collect. We can augment the training data by transforming the examples. This is called **data augmentation**.
- Examples (for visual recognition)
  - translation
  - horizontal or vertical flip
  - rotation
  - smooth warping
  - noise (e.g. flip random pixels)
- Only warp the training, not the test, examples.
- The choice of transformations depends on the task. (E.g. horizontal flip for object recognition, but not handwritten digit recognition.)

## Reducing the Capacity

- Roughly speaking, models with more parameters have more capacity. Therefore, we can try to reduce the number of parameters.
- One approach: reduce the number of layers or the number of parameters per layer.
- Adding a linear **bottleneck layer** is another way to reduce the capacity:



- The first network is strictly more expressive than the second (i.e. it can represent a strictly larger class of functions). (Why?)
- Remember how linear layers don't make a network more expressive? They might still improve generalization.

# Weight Decay

- We've already seen that we can **regularize** a network by penalizing large weight values, thereby encouraging the weights to be small in magnitude.

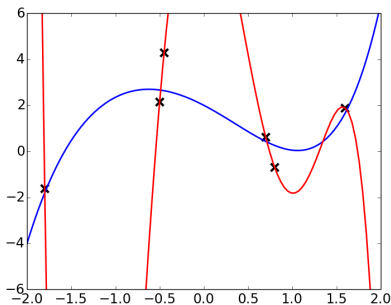
$$\mathcal{E}_{\text{reg}} = \mathcal{E} + \lambda \mathcal{R} = \mathcal{E} + \frac{\lambda}{2} \sum_j w_j^2$$

- We saw that the gradient descent update can be interpreted as **weight decay**:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \left( \frac{\partial \mathcal{E}}{\partial \mathbf{w}} + \lambda \frac{\partial \mathcal{R}}{\partial \mathbf{w}} \right) \\ &= \mathbf{w} - \alpha \left( \frac{\partial \mathcal{E}}{\partial \mathbf{w}} + \lambda \mathbf{w} \right) \\ &= (1 - \alpha \lambda) \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}} \end{aligned}$$

# Weight Decay

Why we want weights to be small:



$$y = 0.1x^5 + 0.2x^4 + 0.75x^3 - x^2 - 2x + 2$$

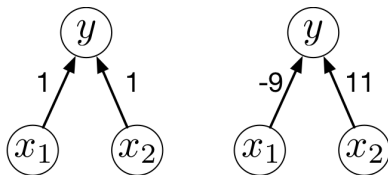
$$y = -7.2x^5 + 10.4x^4 + 24.5x^3 - 37.9x^2 - 3.6x + 12$$

The red polynomial overfits. Notice it has really large coefficients.

# Weight Decay

Why we want weights to be small:

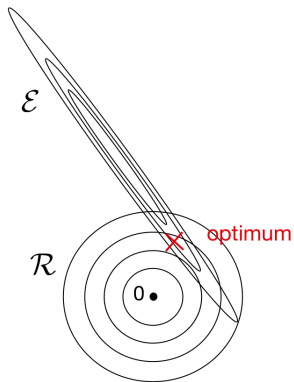
- Suppose inputs  $x_1$  and  $x_2$  are nearly identical. The following two networks make nearly the same predictions:



- But the second network might make weird predictions if the test distribution is slightly different (e.g.  $x_1$  and  $x_2$  match less closely).

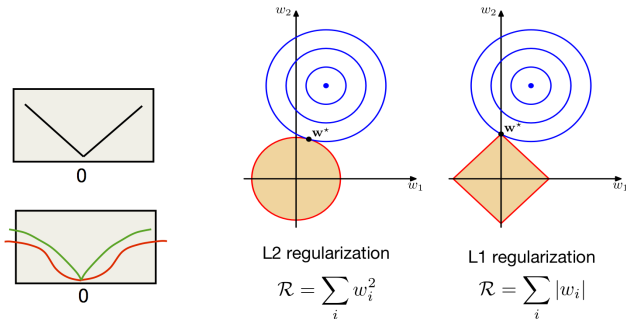
# Weight Decay

- The geometric picture:



# Weight Decay

- There are other kinds of regularizers which encourage weights to be small, e.g. sum of the absolute values.
- These alternative penalties are commonly used in other areas of machine learning, but less commonly for neural nets.
- Regularizers differ by how strongly they prioritize making weights exactly zero, vs. not being very large.



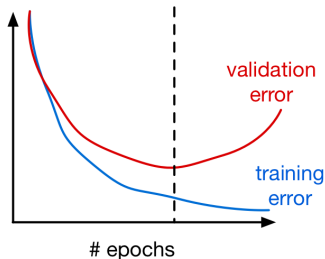
— Hinton, Coursera lectures

— Bishop, *Pattern Recognition and Machine Learning*



# Early Stopping

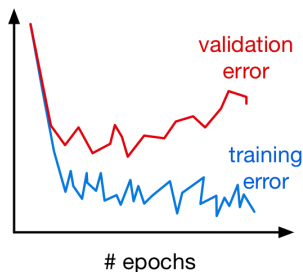
- We don't always want to find a global (or even local) optimum of our cost function. It may be advantageous to stop training early.
- Roughly speaking, training for longer increases the capacity of the model, which might make us overfit.



- **Early stopping:** monitor performance on a validation set, stop training when the validation error starts going up.

# Early Stopping

- A slight catch: validation error fluctuates because of stochasticity in the updates.



- Determining when the validation error has actually leveled off can be tricky.

# Early Stopping

- Why does early stopping work?
  - Weights start out small, so it takes time for them to grow large. Therefore, it has a similar effect to weight decay.
  - If you are using sigmoidal units, and the weights start out small, then the inputs to the activation functions take only a small range of values.
    - Therefore, the network starts out approximately linear, and gradually becomes more nonlinear (and hence more powerful).

# Ensembles

- If a loss function is convex (with respect to the predictions), you have a bunch of predictions, and you don't know which one is best, you are always better off averaging them.

$$\mathcal{L}(\lambda_1 y_1 + \dots + \lambda_N y_N, t) \leq \lambda_1 \mathcal{L}(y_1, t) + \dots + \lambda_N \mathcal{L}(y_N, t) \quad \text{for } \lambda_i \geq 0, \sum_i \lambda_i = 1$$

- This is true no matter where they came from (trained neural net, random guessing, etc.). Note that only the loss function needs to be convex, not the optimization problem.
- Examples: squared error, cross-entropy, hinge loss
- If you have multiple candidate models and don't know which one is the best, maybe you should just average their predictions on the test data. The set of models is called an **ensemble**.
- Averaging often helps even when the loss is nonconvex (e.g. 0–1 loss).

# Ensembles

- Some examples of ensembles:
  - Train networks starting from different random initializations. But this might not give enough diversity to be useful.
  - Train networks on different subsets of the training data. This is called **bagging**.
  - Train networks with different architectures or hyperparameters, or even use other algorithms which aren't neural nets.
- Ensembles can improve generalization quite a bit, and the winning systems for most machine learning benchmarks are ensembles.
- But they are expensive, and the predictions can be hard to interpret.

# Stochastic Regularization

- For a network to overfit, its computations need to be really precise. This suggests regularizing them by injecting noise into the computations, a strategy known as **stochastic regularization**.
- **Dropout** is a stochastic regularizer which randomly deactivates a subset of the units (i.e. sets their activations to zero).

$$h_i = \begin{cases} \phi(z_i) & \text{with probability } 1 - \rho \\ 0 & \text{with probability } \rho, \end{cases}$$

where  $\rho$  is a hyperparameter.

- Equivalently,

$$h_i = m_i \cdot \phi(z_i),$$

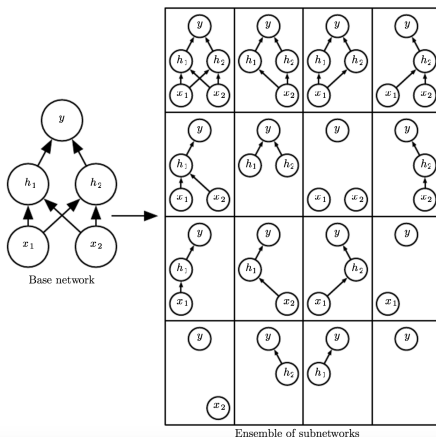
where  $m_i$  is a Bernoulli random variable, independent for each hidden unit.

- Backprop rule:

$$\bar{z}_i = \bar{h}_i \cdot m_i \cdot \phi'(z_i)$$

# Stochastic Regularization

- Dropout can be seen as training an ensemble of  $2^D$  different architectures with shared weights (where  $D$  is the number of units):



— Goodfellow et al., *Deep Learning*

# Stochastic Regularization

- Dropout can help performance quite a bit, even if you're already using weight decay.
- Lots of other stochastic regularizers have been proposed:
  - DropConnect drops connections instead of activations.
  - Batch normalization (mentioned last week for its optimization benefits) also introduces stochasticity, thereby acting as a regularizer.
  - The stochasticity in SGD updates has been observed to act as a regularizer, helping generalization.
    - Increasing the mini-batch size may improve training error at the expense of test error!



# Our Bag of Tricks

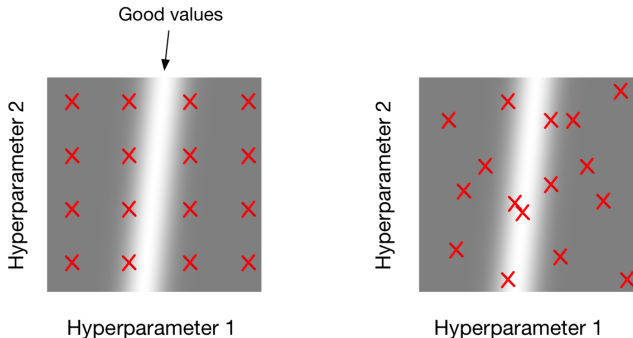
- Techniques we just covered:
  - data augmentation
  - reduce the capacity (e.g. number of parameters)
  - weight decay
  - early stopping
  - ensembles (combine predictions of different models)
  - stochastic regularization (e.g. dropout, batch normalization)
- The best-performing models on most benchmarks use some or all of these tricks.
- Many of these techniques involve hyperparameters. How do we choose these?

# Choosing Hyperparameters

- Many hyperparameters are relevant to generalization
  - number of layers
  - number of units per layer
  - $L_2$  weight cost
  - dropout probability
- Ideally, we want to choose hyperparameters which perform the best on a validation set.
- Ways to approximate this
  - manual (“graduate student descent”)
  - grid search
  - random search
  - Bayesian optimization (covered later in the course)
  - use Autograd to differentiate through the whole training procedure, and do gradient descent on hyperparameters (only practical for very small problems)

# Choosing Hyperparameters

- Random search can be more efficient than grid search when some of the hyperparameters are unimportant:



- But grid search can be more reproducible.