# Model Capacity, Regularization, and Hyperparameters

**Danna Gurari**

University of Colorado Boulder

Spring 2025

# Review

- Last lecture:
  - Motivation: effective gradients for learning
  - Initializing parameters
  - Initializing data
  - Following the gradient (optimization)
  - Programming tutorial

- Assignments (Canvas):
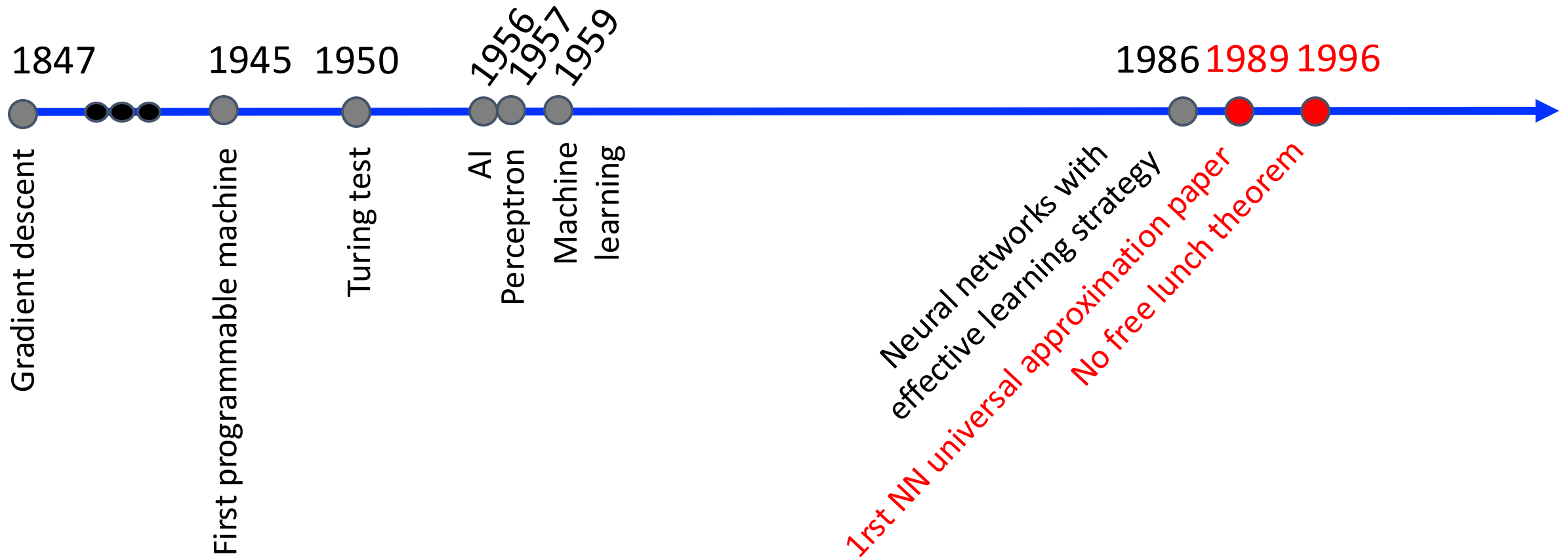  - Problem set 2 due on Tuesday

- Questions?

# Today's Topics

- Model capacity: how it affects learning

- Regularization: learning methods for improving model generalization

- Hyperparameter selection: tuning to improve model performance

- Programming tutorial

# Today's Topics

- Model capacity: how it affects learning

- Regularization: learning methods for improving model generalization

- Hyperparameter selection: tuning to improve model performance

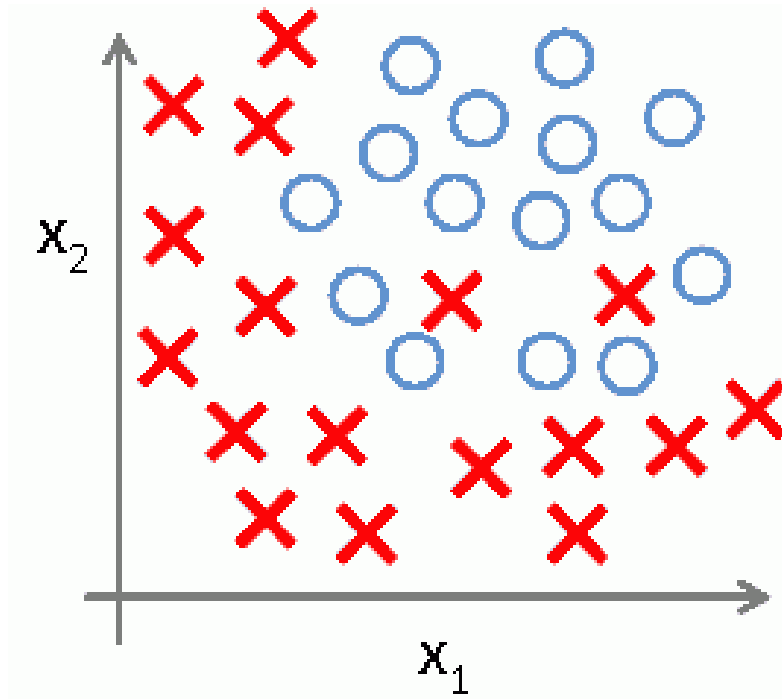- Programming tutorial

# Historical Context: Motivating Theory



1847 — Gradient descent
1945 — First programmable machine
1950 — Turing test
1956 — AI
1957 — Perceptron
1959 — Machine learning
1986 — Neural networks with effective learning strategy
1989 — 1rst NN universal approximation paper
1996 — No free lunch theorem

Hornik, Stinchcombe and White. Multilayer feedforward networks are universal approximators. Neural Networks, 1989

# What model design should we use?

"no free lunch theorem… no machine learning algorithm is universally any better than any other." - Ch. 5.2.1 of Goodfellow book

"The universal approximation theorem means that regardless of what function we are trying to learn, we know that a large MLP [multilayer perceptron] will be able to *represent* this function." - Ch. 6.4.1 of Goodfellow book

# Model Capacity: Recall Class Exercise
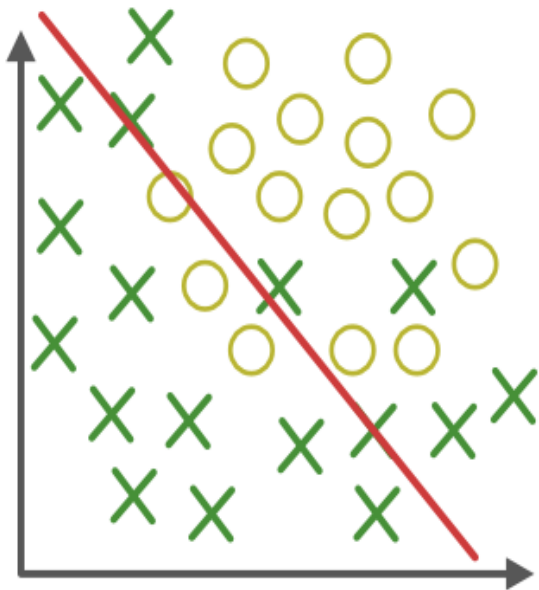
- Model-based approach: separate x from o

Class volunteer:
1) Draw a straight line (linear equation)
2) Draw a parabola (quadratic equation)
3) Draw any curve

Models with increasing representational capacity

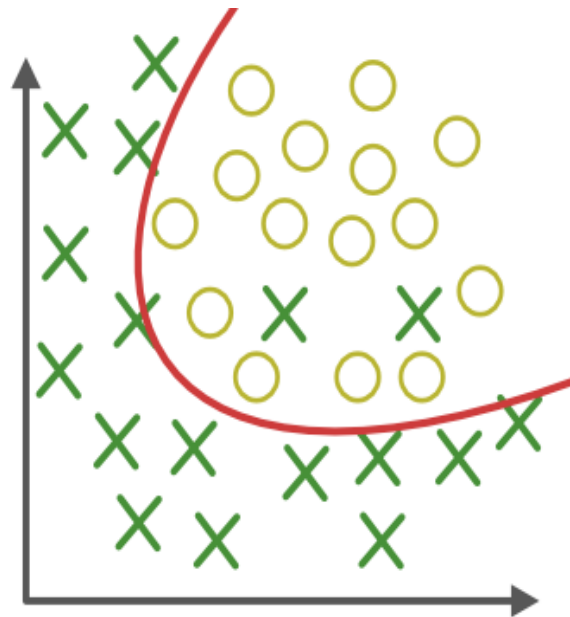# Model Capacity
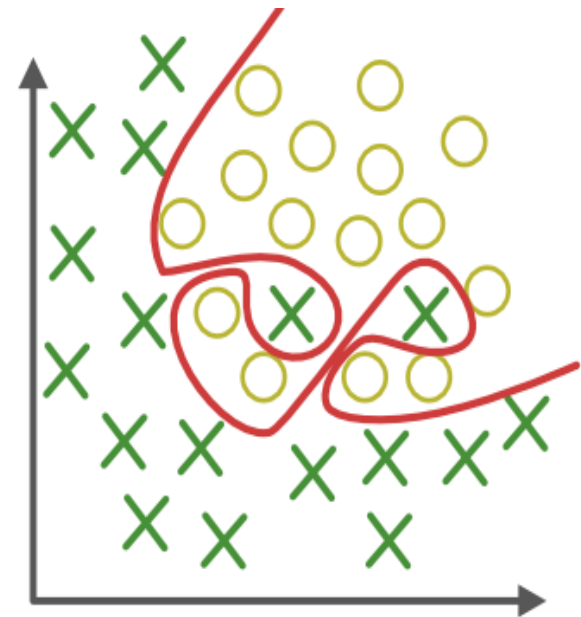
Which model would you choose to separate x from o?

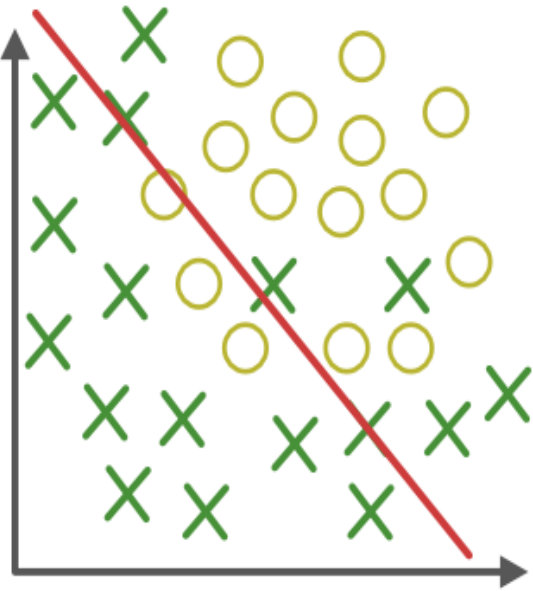(a)                                    (b)                                    (c)
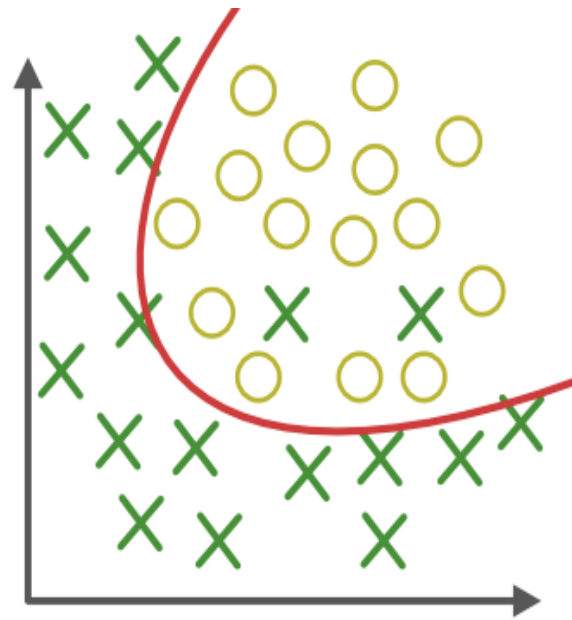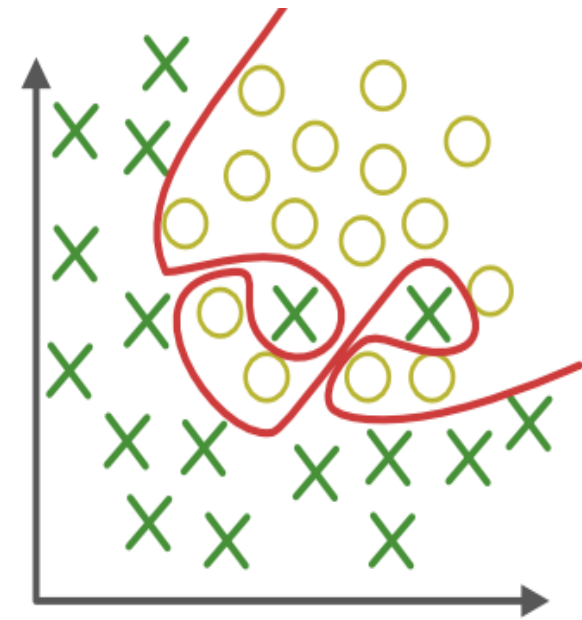
# Model Capacity

Underfits: too simple
to explain the data
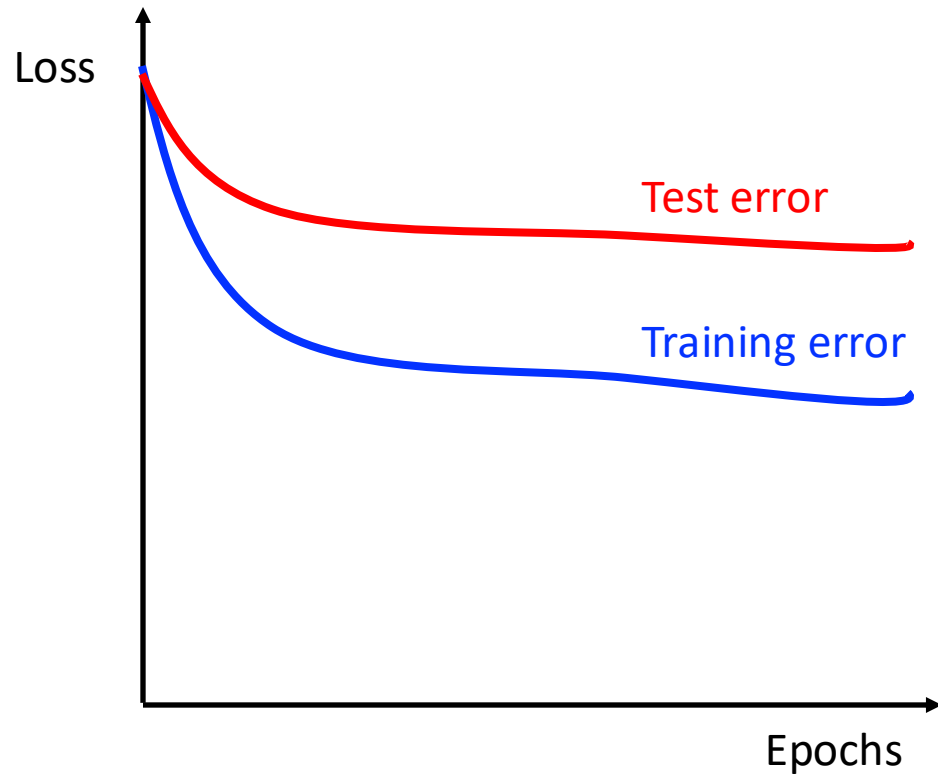
(a)

(b)

Overfits: too complex to
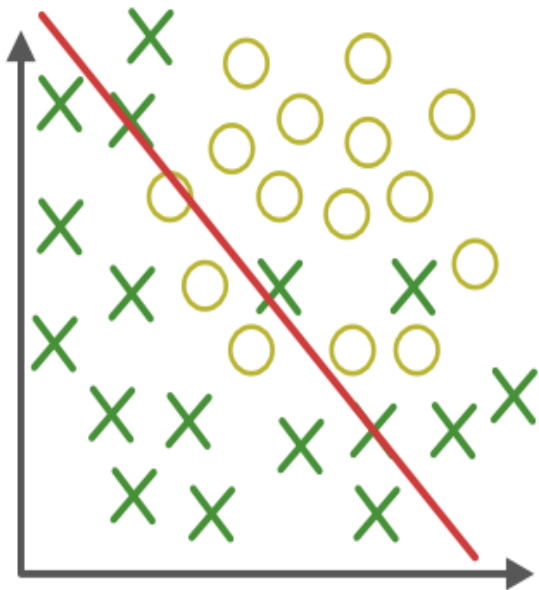generalize to a test set

(c)

# Model Capacity: Detecting Underfitting



Analyze **learning curves** for models tested on training data

- What happens to training data error as number of training steps increases?
  - Error remains high
- What happens to test data error as number of training steps increases?
  - Error remains high
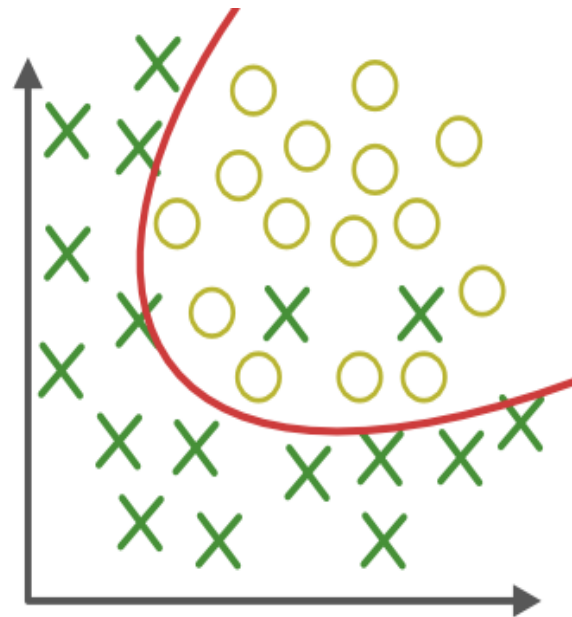
# Model Capacity

Underfits: too simple
to explain the data
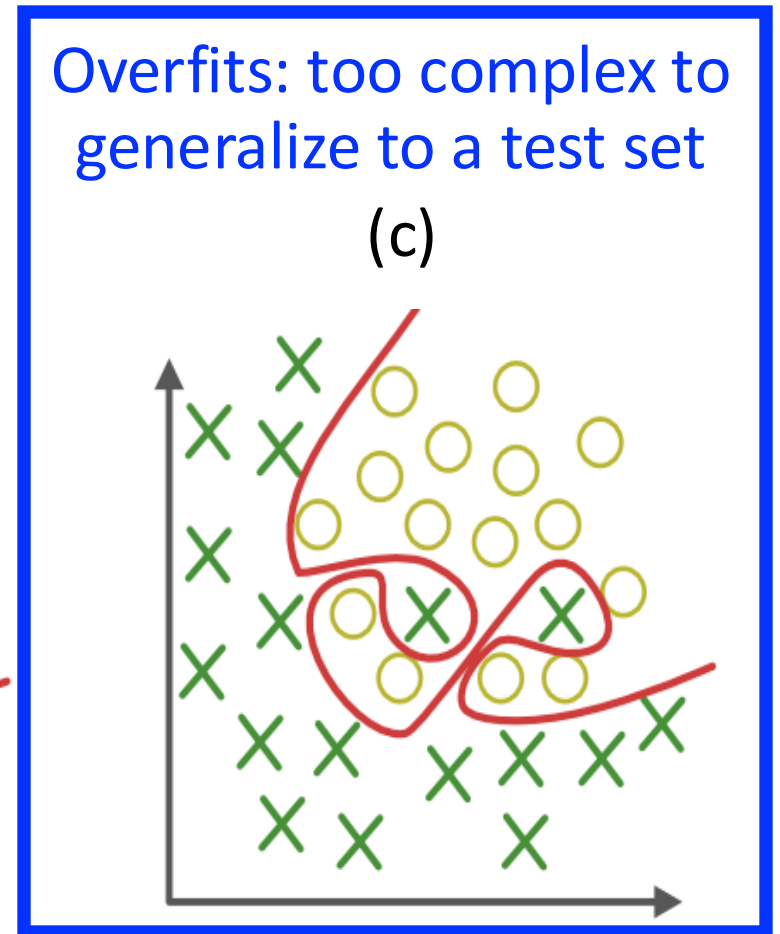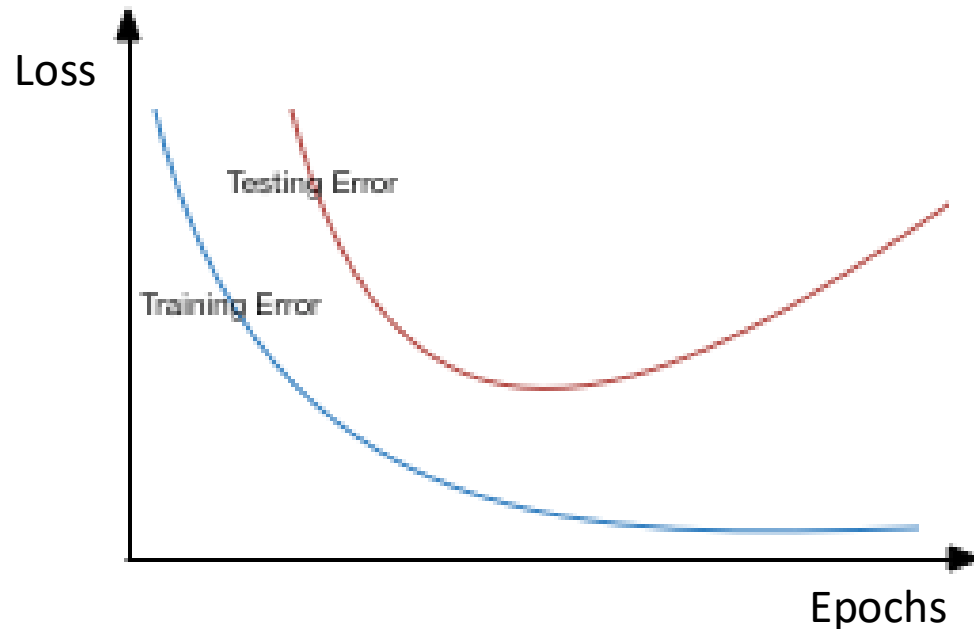
(a)

(b)

Overfits: too complex to
generalize to a test set

(c)

# Model Capacity: Detecting Overfitting



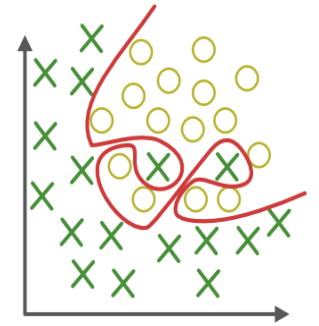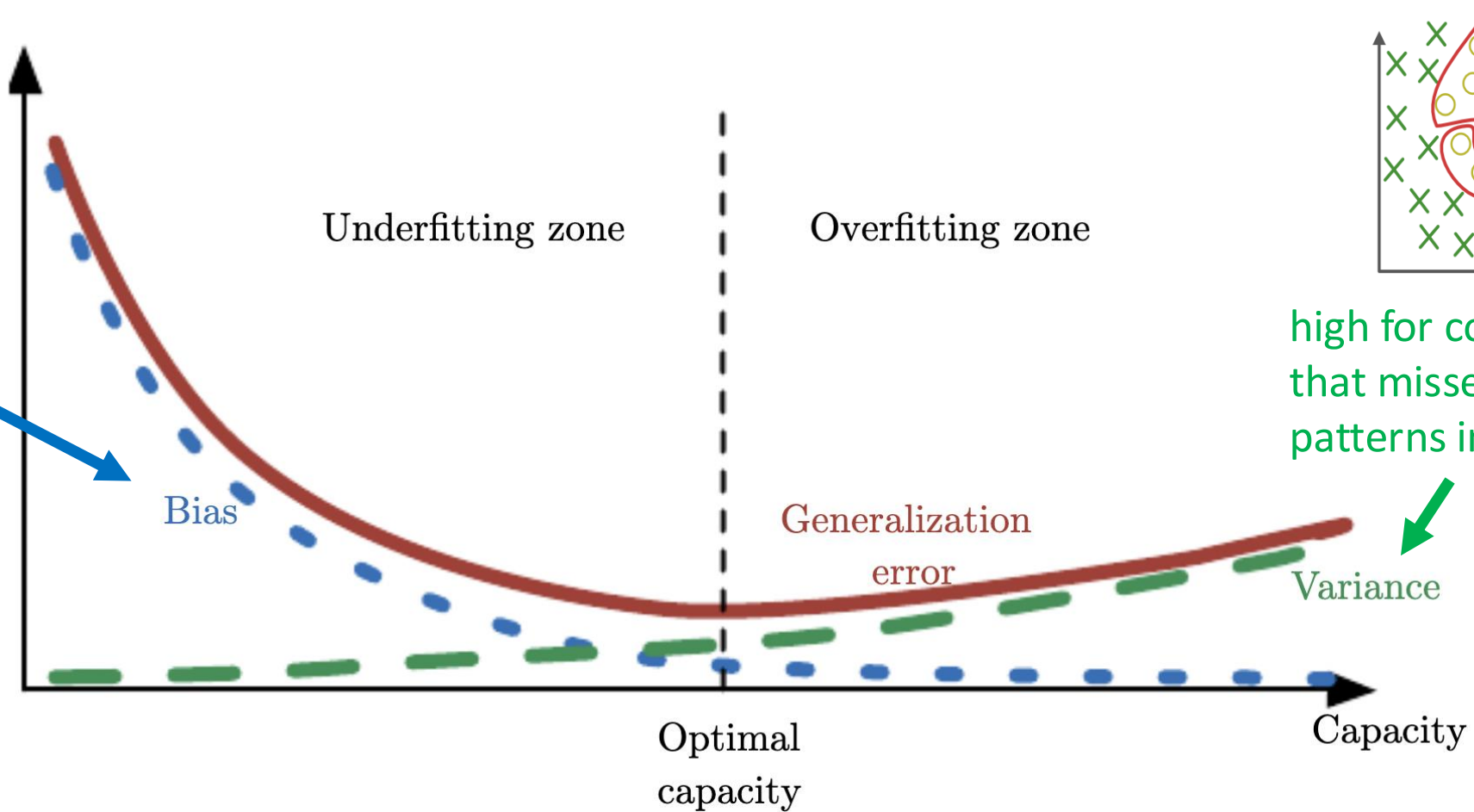Analyze **learning curves** for model on <span style="color:blue">training data</span> and <span style="color:red">test data</span>

- What happens to <span style="color:blue">training data</span> error as number of training steps increases?
    - Error shrinks
- What happens to <span style="color:red">test data</span> error as number of training steps increases?
    - Error shrinks and then grows
- Why does <span style="color:blue">training error</span> ***shrink*** and <span style="color:red">test error</span> ***grow***?
    - Models ***noise*** in the training data (i.e., "overfitting") at the expense of knowledge that generalizes
- What can cause noise in a dataset?
    - e.g., incorrect data entry/labeling, hardware measurement error (caution: some outliers are data points we want models to learn)

# Model Capacity: Overfitting vs Underfitting



high for simple model that assumes simpler patterns in the data (not related to "bias" parameter):

high for complex model that misses simpler patterns in the data

Underfitting zone

Overfitting zone

Bias

Generalization error

Variance

Optimal capacity

Capacity

Often discussed with respect to a **bias-variance** trade-off

# Model Capacity: Overfitting vs Underfitting



Can shift between underfitting and overfitting by adding/removing parameters

# Model Capacity: Overfitting vs Underfitting



Often, we err on over-parameterized models capable of overfitting and then regularize them

Ian Goodfellow, Yoshua Bengio, and Aaron Courville; Deep Learning, 2016

# During Training, You Should Ask Yourself: What Does the Observed Performance Mean?

- Loss and performance curves can signal how well training is going

# Today's Topics

- Model capacity: how it affects learning

- **Regularization: learning methods for improving model generalization**

- Hyperparameter selection: tuning to improve model performance

- Programming tutorial

# What is Regularization?

- "any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error."- Ch. 5.2 of Goodfellow book

Regularization techniques will be a focus for much of the course

# e.g., Regularize Over-Parameterized Model

(a)

(b)

(c)

# Example 1: Early Stopping During Training



Error

Testing Error

Training Error

Use This Model

Training steps

# Example 2: Parameter Norm Penalty

Smooth a model's decision boundaries by incentivizing against large weights

(models can learn to represent noise with correlated large positive and negative weights that usually just cancel each other out)

# Example 2: Parameter Norm Penalty

Add penalty term to objective function; e.g., when *minimizing* **s**um of **s**quared **e**rrors

- L2 norm: penalize squared weight values

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \alpha\boxed{\sum_{j=1}^{m} w_j^2}$$

- L1 norm: penalize absolute weight values

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \alpha\boxed{\sum_{j=1}^{m} |w_j|}$$

Analogy: Belt for Big Pants



- *Note: penalizes only weight (not biases) and can apply per layer or globally*

# Example 2: Parameter Norm Penalty

Add penalty term to objective function; e.g., when *minimizing* sum of squared errors

- L2 norm: penalize squared weight values

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \boxed{\alpha}\boxed{\sum_{j=1}^{m} w_j^2}$$

- L1 norm: penalize absolute weight values

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \boxed{\alpha}\boxed{\sum_{j=1}^{m} |w_j|}$$

Analogy: Belt for Big Pants



- Hyperparameter determines contribution of norm penalty term (e.g., belt tightness)

# Example 2: Parameter Norm Penalty

Add penalty term to objective function; e.g., when *minimizing* sum of squared errors

- L2 norm: penalize squared weight values

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \boxed{\alpha}\boxed{\sum_{j=1}^{m} w_j^2}$$

- L1 norm: penalize absolute weight values

Intuitively, larger alpha values prioritizes having weights closer to 0 instead of minimizing sum of squared errors

$$Error = \boxed{\sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2} + \boxed{\alpha}\boxed{\sum_{j=1}^{m} |w_j|}$$

- Hyperparameter determines contribution of norm penalty term (e.g., belt tightness)

# Example 2: Parameter Norm Penalty

Add penalty term to objective function; e.g., when *minimizing* sum of squared errors

- L2 norm: penalize squared weight values

$$Error = \sum_{i=1}^{n} (y^{(i)} - \widehat{y}^{(i)})^2 + \alpha \boxed{\sum_{j=1}^{m} w_j^2}$$

- L1 norm: penalize absolute weight values

$$Error = \sum_{i=1}^{n} (y^{(i)} - \widehat{y}^{(i)})^2 + \alpha \boxed{\sum_{j=1}^{m} |w_j|}$$

- Only change for learning is need gradient for norm penalty (shown in assigned reading)
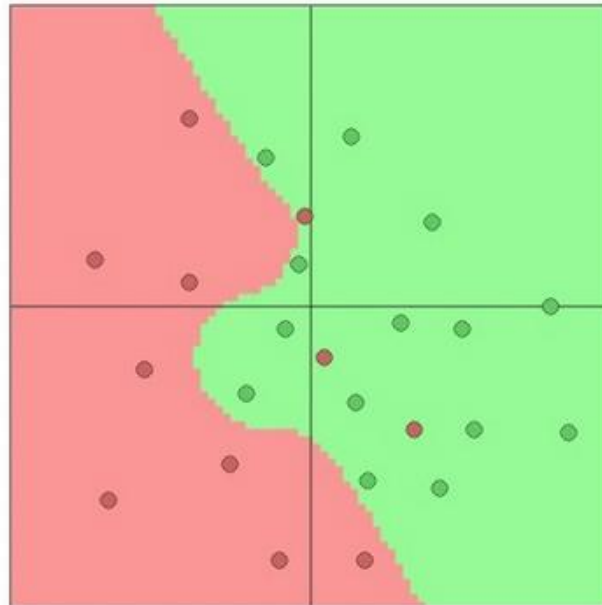
# Example 2: Parameter Norm Penalty

Which model had the largest value for alpha (i.e., norm penalty contribution), given that this value was the only difference when training the models?

$$Error = \sum_{i=1}^{n}(y^{(i)} - \widehat{y}^{(i)})^2 + \boxed{\alpha}\sum_{j=1}^{m} w_j^2$$

(a)                               (b)                               (c)

# Example 2: Parameter Norm Penalty (Geometric Interpretation in 2D)



Minimizes sum of squared errors cost

Note: L2 commonly used in practice

L1

L2

$\beta_2$

$\hat{\beta}$

$\beta_2$

$\hat{\beta}$

Minimizes cost + penalty

Contour of least square error function

$\beta_1$

$\beta_1$

Contour of regularization constraint functions shown below both plots

Minimizes penalty term

$$|\beta_1| + |\beta_2| \le t$$

$$\beta_1^2 + \beta_2^2 \le t^2$$

# Early Stopping vs Parameter Norm Penalties

Similar behavior when model parameters initialized around origin; e.g.,

Early Stopping　　　　　Parameter Norm Penalty (L2)

# Today's Topics

- Model capacity: how it affects learning

- Regularization: learning methods for improving model generalization

- **Hyperparameter selection: tuning to improve model performance**

- Programming tutorial

# Model Design Decisions

**Hyperparameters (selected); e.g.,**
- Number of layers
- Number of units in each layer
- Activation function
- Batch size
- Learning rate
- …

**Parameters (learned)**
- Weights
- Biases
- …

**Challenge:** how to design a model without repeatedly observing the test data (which leads to overfitting)?
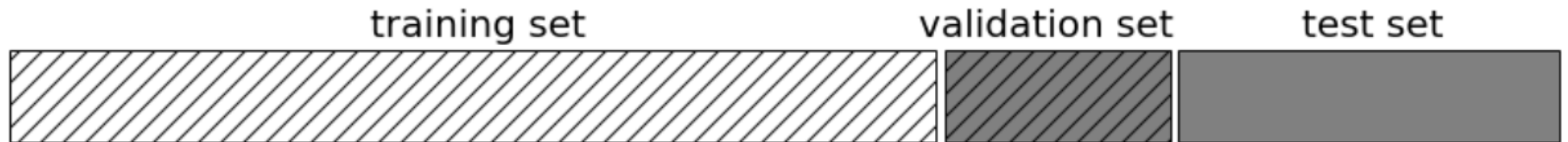
# Recall: Our Goal is to Design Models that **Generalize** Well to New, Previously Unseen Examples (Test Data)



**Challenge:** how to design a model without repeatedly observing the test data (which leads to overfitting)?

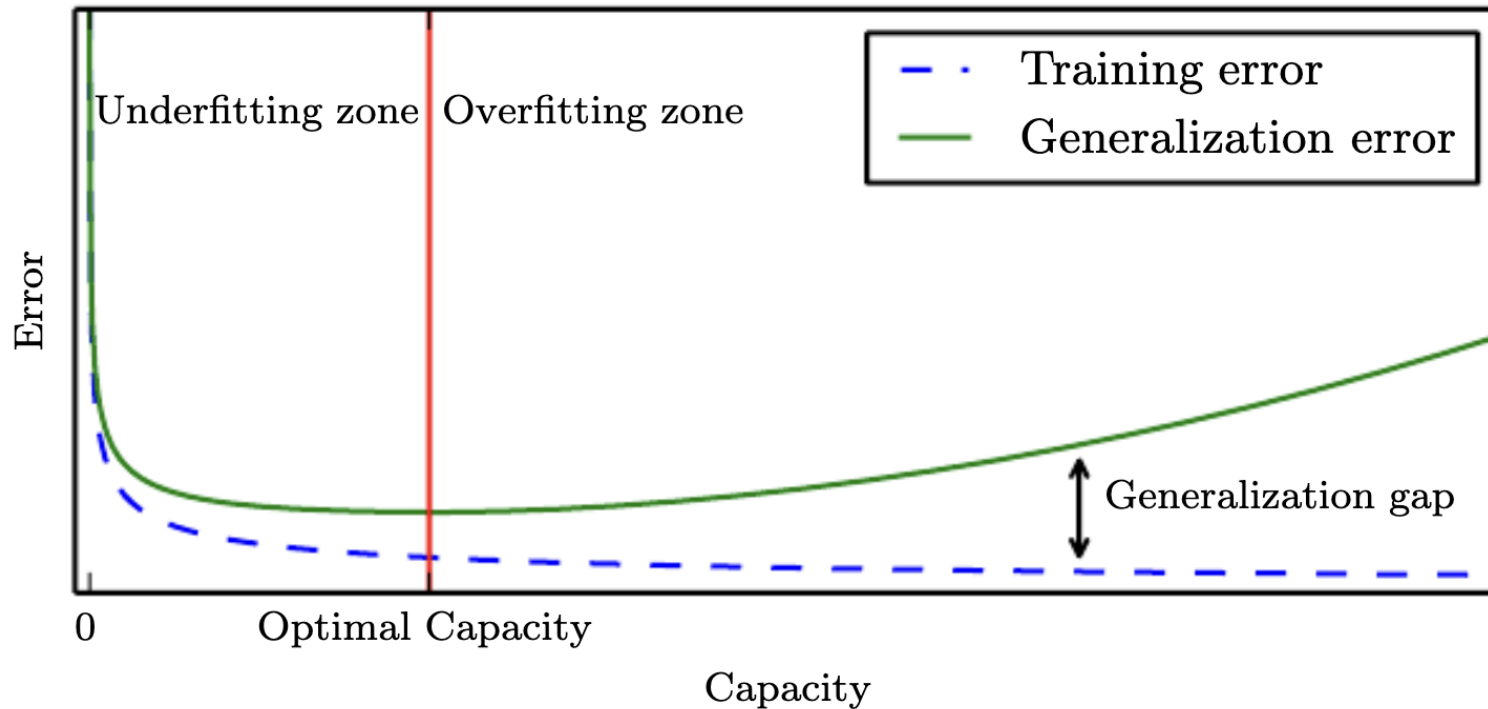# Validation Dataset for Hyperparameter Tuning

- Dataset divided into three splits; e.g., using a 60%/20%/20% train/val/test split



- **Hyperparameter selection**: test on validation set to identify best hyperparameters
- **Final model**: train on training AND validation splits with best hyperparameters
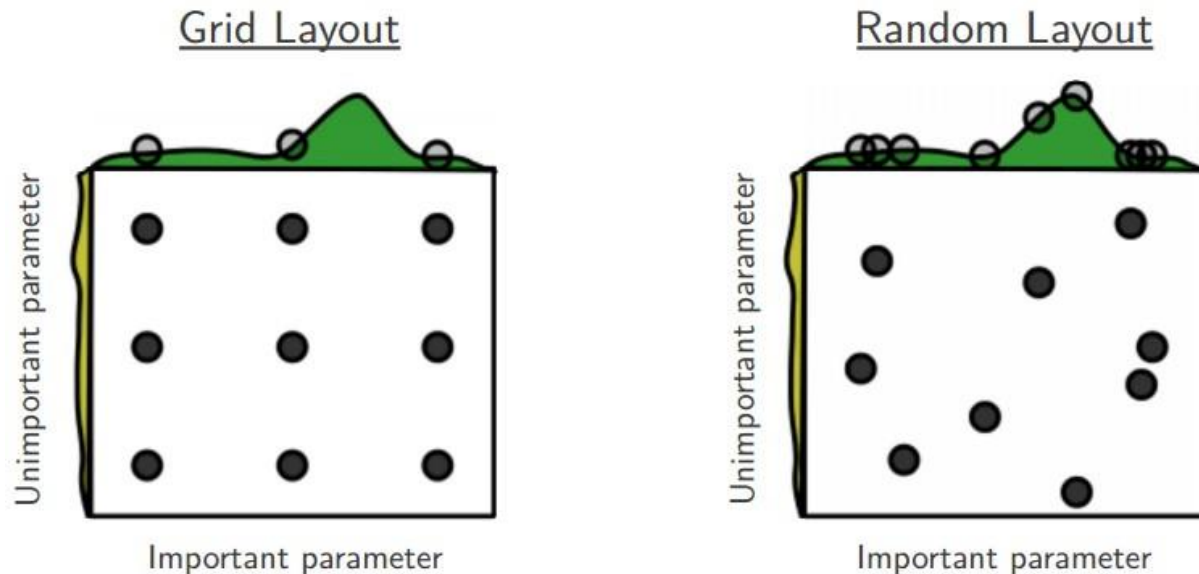
# Hyperparameter Tuning for Optimal Capacity



Tweaking a hyperparameter can shift model between underfitting and overfitting; e.g., impact of increasing vs decreasing…

- # of hidden layers and hidden units (model architecture)?
- weight decay coefficient (regularization strength)?
- learning rate (model training)?

Ian Goodfellow, Yoshua Bengio, and Aaron Courville; Deep Learning, 2016

# Hyperparameter Tuning Approaches

- **Automatic**: extensive and so computationally costly; e.g.,



Grid Layout / Random Layout (Unimportant parameter vs Important parameter)

As exemplified, grid search is inferior when selecting multiple hyperparameters since **fewer values are tested** for each (in this example, two hyperparameters); random search better supports identifying impactful hyperparameters

- **Manual**: leverages "art" of knowing how to effectively train
  - understanding relationship between hyperparameters, training error, generalization error, and computational constraints (e.g., GPUs, memory)

Bergstra et al; Random Search for Hyper-Parameter Optimization, 2012

# Discussion: What Do You Think Is Happening and What Might Be a Fix?

Error on training set is larger than your target error rate

# Today's Topics

• Model capacity: how it affects learning

• Regularization: learning methods for improving model generalization

• Hyperparameter selection: tuning to improve model performance

• **Programming tutorial**

# Today's Topics

- Model capacity: how it affects learning

- Regularization: learning methods for improving model generalization

- Hyperparameter selection: tuning to improve model performance

- Programming tutorial