# 6.390 Midterm Review

Spring 2025 Haley Nakamura

## Topic 1: IntroML

- ML: making predictions from data
- Training / validation / testing data
- Supervised learning
- Hypothesis / hypothesis class
- Learning algorithm
- Evaluation metrics
  - Loss functions
  - Overfitting vs. underfitting

$$egin{split} \mathcal{D}_{ ext{train}} &= \left\{ \left(x^{(1)},y^{(1)}
ight),\ldots,\left(x^{(n)},y^{(n)}
ight)
ight\} \ & x o ight[h] o y \ \mathcal{D}_{ ext{train}} &\longrightarrow igl[ ext{learning alg}\left(\mathcal{H}
ight)igg] o h \ \mathcal{E}_{ ext{train}}(h;\Theta) &= rac{1}{n}\sum_{i=1}^n \mathcal{L}(h(x^{(i)};\Theta),y^{(i)}) \end{split}$$

#### Topic 2: Linear Regression

- Label: continuous real number

$$y=h(x; heta, heta_0)= heta^Tx+ heta_0$$

- Linear hypothesis class
- Objective function (e.g. MSE)

$$f( heta, heta_0) = rac{1}{n}\sum_{i=1}^n \left( heta^T x^{(i)} + heta_0 - y^{(i)}
ight)^2$$

n

- Closed-form/analytical solution (optimal)
- What if  $X^T X$  is not invertible?

$$heta^* = ig(X^T Xig)^{-1} X^T Y$$

- Objective function has "half-pipe" shape instead of "bowl" shape

J

- Cases:
  - More features than data points
  - Features are linearly dependent

#### Topic 2: Linear Regression

- Add regularizer term

$$J_{ ext{ridge}}( heta, heta_0) = rac{1}{n}\sum_{i=1}^n \left( heta^T x^{(i)} + heta_0 - y^{(i)}
ight)^2 + \lambda \| heta\|^2$$

-  $\lambda$  as a hyperparameter

- Generalizable

$$heta_{ ext{ridge}} = ig(X^TX + n\lambda Iig)^{-1}X^TY$$

- Cross-validation

#### Topic 3: Gradient Descent

- Gradient vector (analytically and conceptually)
- Algorithm and update step  $\theta^{(t+1)} = \theta^{(t)} \eta \nabla_{\theta} J(\theta)|_{\theta = \theta^{(t)}}$
- $\eta$  as a hyperparameter
- Termination criterion
- Smooth, convex, sufficiently small learning rate, enough iterations, global min exists?
  - Converge to global minimum!
- What if these assumptions are violated?
- Stochastic gradient descent (what changes?)

## Topic 4: Classification

 $h(x; heta, heta_0) = ext{step}( heta^T x + heta_0)$ 

 $z= heta^Tx+ heta_{ extsf{n}}$ 

- Binary classification: {1, 0}
- Linear classifier
  - Normal vector
- Logistic classifier
  - Sigmoid, NLL Loss
  - Magnitude of parameters?
- Multiclass:
  - Pick one: softmax, NLLM Loss
  - Pick many: multiple sigmoids

## $h(x; heta, heta_0) = ext{step}( heta^T x + heta_0)$

## Topic 4: Classification

- Binary classification: {1, 0}
- Linear classifier
  - Normal vector
- Logistic classifier
  - Sigmoid, NLL Loss
  - Magnitude of parameters?
- Multiclass:
  - Pick one: softmax, NLLM Loss
  - Pick many: multiple sigmoids



## Topic 4: Classification

- Binary classification: {1, 0}
- Linear classifier
  - Normal vector
- Logistic classifier
  - Sigmoid, NLL Loss
  - Magnitude of parameters?
- Multiclass:
  - Pick one: softmax, NLLM Loss
  - Pick many: multiple sigmoids



## Topic 4: Classification

- Binary classification: {1, 0}
- Linear classifier
  - Normal vector
- Logistic classifier
  - Sigmoid, NLL Loss
  - Magnitude of parameters?
- Multiclass:
  - Pick one: softmax, NLLM Loss
  - Pick many: multiple sigmoids



#### Topic 5: Features

- Nonlinear feature transformations
- (k<sup>th</sup> order) Polynomial basis

#### Not linearly separable in x space



Linearly separable in  $\phi(x) = x^2$  space

- Computation graph
  - Neurons  $\rightarrow$  Layers
    - $\rightarrow$  Network



- Computation graph
  - Neurons  $\rightarrow$  Layers
    - $\rightarrow$  Network



- Computation graph -
  - Neurons  $\rightarrow$  Layers -
    - $\rightarrow$  Network
- Network as a hypothesis -

(forward-pass)



- Output layer design
  - Dimension (# output neurons), activation function, loss function
- Hand-crafting weights to match a function

- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



 $\partial \mathcal{L}(g,y)$ 

- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)  $\frac{\partial \mathcal{L}(g, y)}{\partial Z^2}$



- Backward-pass (backpropagation)
  - Shared partial derivatives when evaluating gradients (w.r.t. different network parameters)



 $\partial \mathcal{L}(g,y)$ 

 $\partial Z^2$ 

