

Intro to Deep Learning

Lecture 2:
Machine Learning Fundamentals

Before getting into deep
learning...

Before getting into deep
learning...

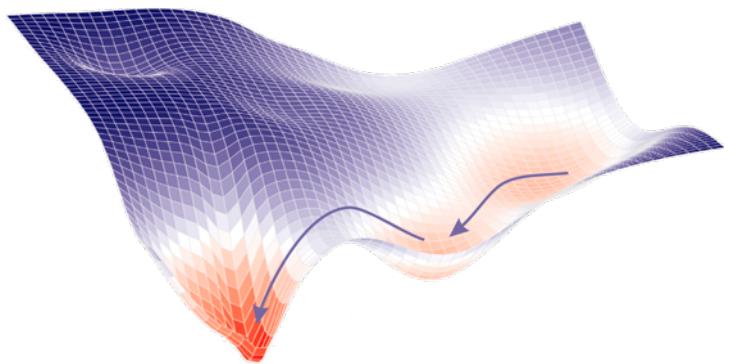
Let us recall ML basics...

Course Topics

Machine Learning



Optimization



Example: image classification (multi-class)

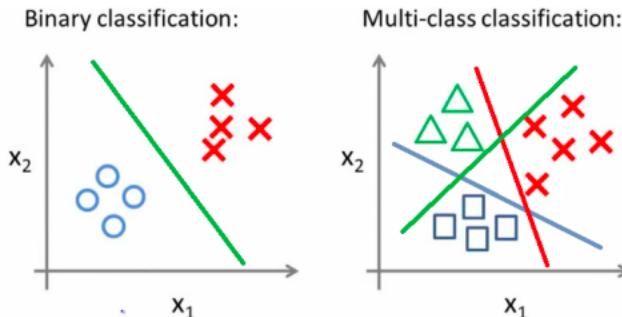


~15 million data points across ~22,000 classes

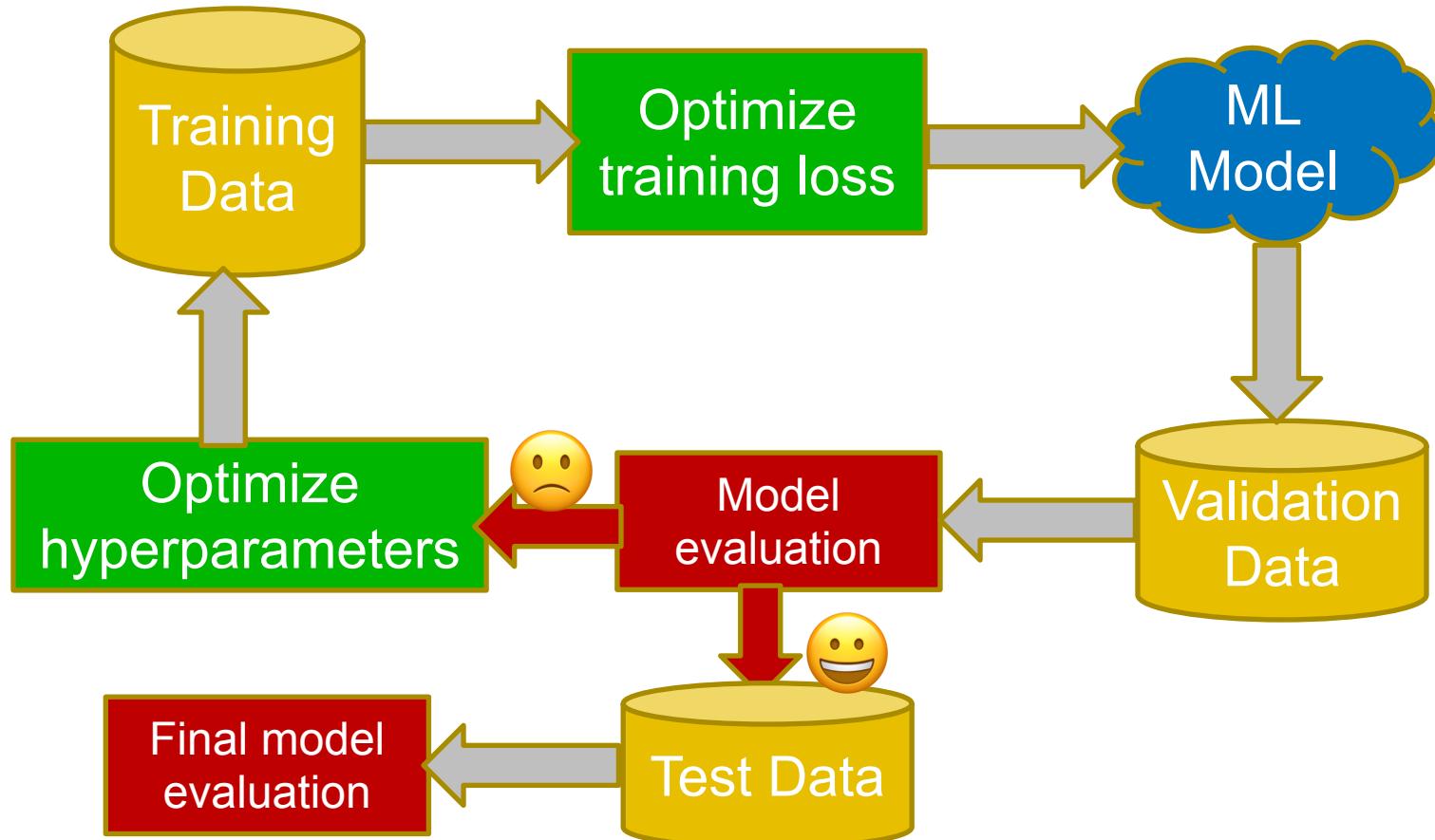
ImageNet figure from vision.stanford.edu

A few terminologies

- **Training data:** images given for learning
- **Validation/Test data:** images to evaluate models
- **Binary classification:** classify into two classes
- **Muticlass classification:** classify into >2 classes

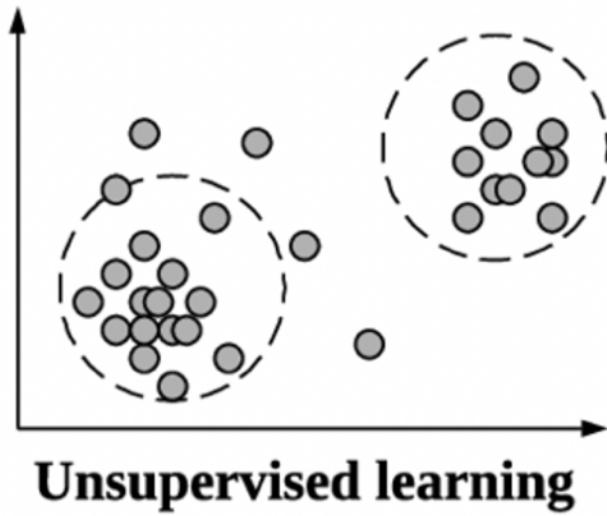
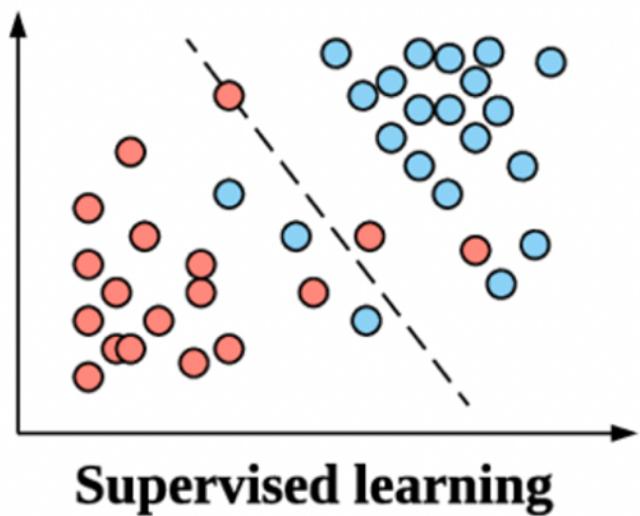


Machine Learning Pipeline



Ingredient 1: Data

Data



Data

Labels

airplane



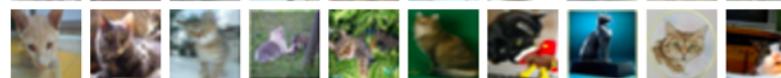
automobile



bird



cat



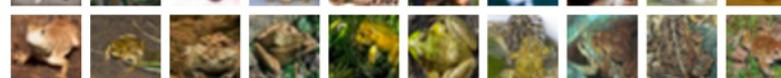
deer



dog



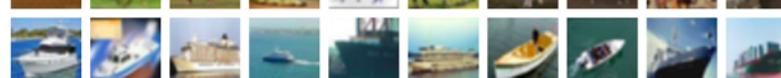
frog



horse



ship



truck

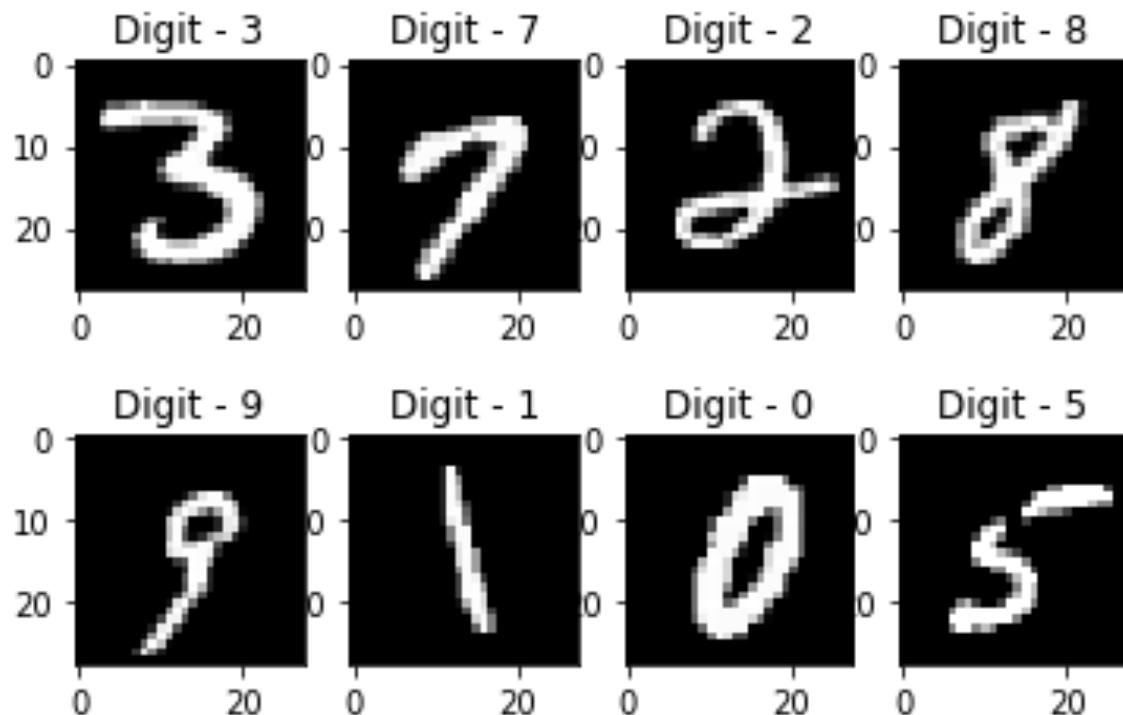


Inputs

Inputs are
 32×32 color images

Notation: A dataset of size n is denoted as $(x_i, y_i)_{i=1}^n$

MNIST data



MNIST data

Input dimension
 $28 \times 28 = 784$

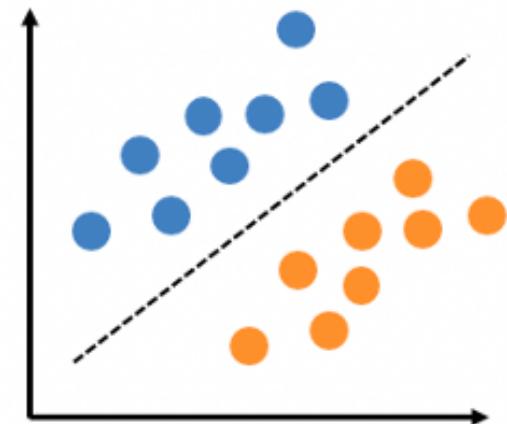
Machine Learning Model

An **ML Model** is a function that outputs a prediction for a given input.

Ex: Linear classifier

Let \hat{y}_f be the label output by f

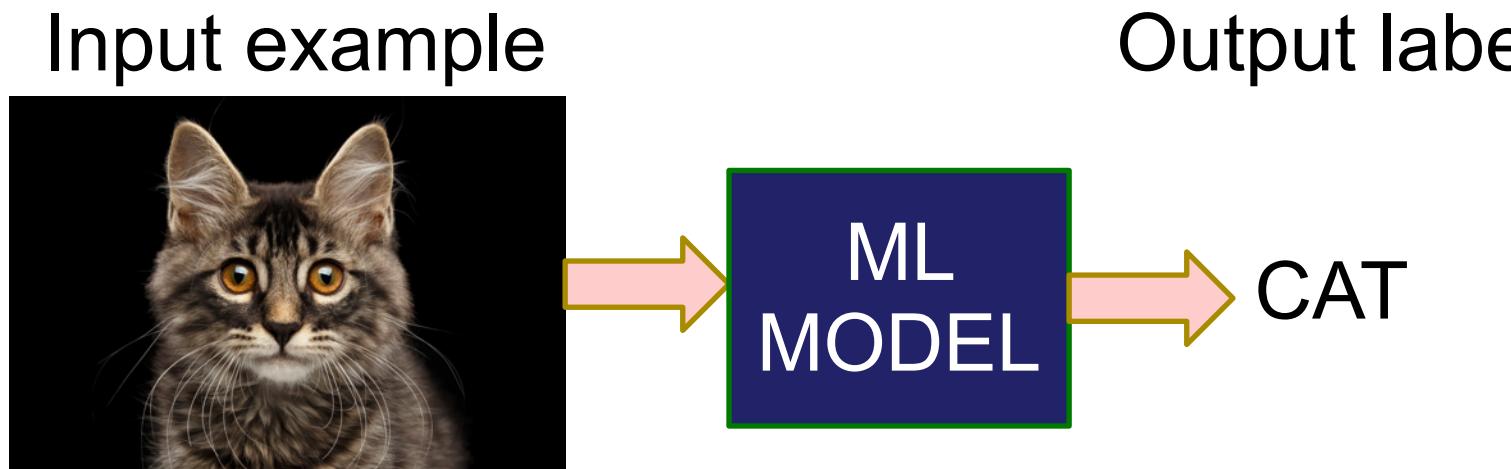
- $\hat{y}_f = 1$ if $f(x) \geq 0$
- $\hat{y}_f = -1$ else



In other words, $\hat{y}_f = \text{sign}(f(x))$

Goal (in supervised ML): "Build a model that can predict the label of the input examples."

Goal (in supervised ML): "Build a model that can predict the label of the input examples."



Goal (in supervised ML): "Build a model that can predict the label of the input examples."

- Given training data $(x_i, y_i)_{i=1}^n$
- Find $f: x \mapsto y$ using training data
- Such that f is (approximately) correct on **unseen** test data

Ingredient 2: Machine Learning Model

An **ML Model** is a function f that outputs a prediction for a given input.

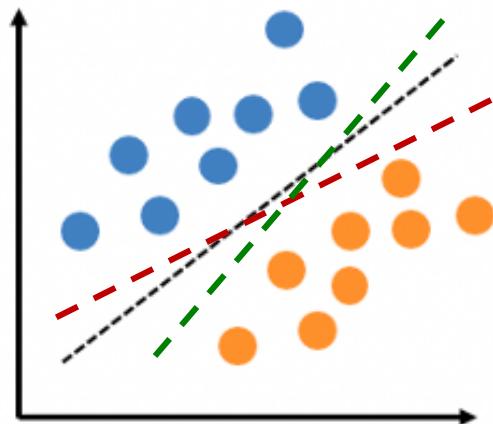
Ex: Input example: $x \in \mathbb{R}^d$

A linear model with weights $w \in \mathbb{R}^d$ outputs

$$f_\theta(x) := \langle w, x \rangle = w^\top x \in \mathbb{R}$$

Ingredient 3: Loss function

There are many ML models



Q: How do you find a good ML model?

A: Minimize a loss function ℓ

Loss function

ℓ takes a prediction $f(x)$ and a label y as input

- **Ex: Squared loss**

$$\ell(f(x), y) = (f(x) - y)^2$$

- **Ex: Cross-entropy loss (multi-class)**

$$\ell(f(x), y) = - \sum_{i \in [C]} y_i \log(p_i), \quad \text{where} \quad p_i = \frac{e^{f_i(x)}}{\sum_{j=1}^n e^{f_j(x)}}$$

- **Ex: 0/1 loss (0-1 loss)**

$$\ell(f(x), y) = \mathbb{I}(\text{sign}(f(x)) \neq y)$$

↗ 1 if $y \neq \text{sign}(f(x))$
↘ 0 if $y = \text{sign}(f(x))$

More examples to come...

Finding the best model

- **Statistical learning:** Data is generated from a distribution

$$(x, y) \sim \mathcal{D}$$

- The model f has the performance

$$\mathcal{L}(f) = \mathbb{E}_{\mathcal{D}}[\ell(f(x), y)]$$

Observation: $\mathbb{E}_{\mathcal{D}}[\mathbb{I}(y \neq \text{sign}(f(x)))] = \mathbb{P}(y \neq \text{sign}(f(x)))$

- We can find best model from a **hypothesis set** \mathcal{F} via

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f)$$

Hypothesis sets

- The model f has the performance

$$\mathcal{L}(f) = \mathbb{E}_{\mathcal{D}}[\ell(f(x), y)]$$

Observation: $\mathbb{E}_{\mathcal{D}}[\mathbb{I}(y \neq \text{sign}(f(x)))] = \mathbb{P}(y \neq \text{sign}(f(x)))$

- We can find best model from a hypothesis set \mathcal{F} via

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f)$$

- Example hypothesis sets

- Linear classifiers
- Neural networks
- Decision trees
- ...

Training Loss

Our aim is typically finding an accurate model i.e.

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f) \quad \text{where} \quad \mathcal{L}(f) = \mathbb{E}_{\mathcal{D}}[\ell(f(x), y)]$$

- **Challenge 1:** We don't have \mathcal{D} . Instead, we only have training data sampled from \mathcal{D}

$$(x_i, y_i)_{i=1}^n \sim \mathcal{D}$$

- **Solution:** Minimize the training loss to approximate $\mathcal{L}(f)$

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \widehat{\mathcal{L}}(f) \quad \text{where} \quad \widehat{\mathcal{L}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

Hopefully \hat{f} will be a good proxy for f^* ! Rigorous proof via concentration inequalities!

Machine Learning Pipeline

1. Data: Collect training dataset
2. Model: Fix hypothesis set \mathcal{F} and loss ℓ
3. Optimize: Minimize the training loss

Example ML Pipeline

1. Data: Collect $(x_i, y_i)_{i=1}^n$



2. Model: Choose a neural network architecture and pick cross-entropy loss
3. Optimize

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \widehat{\mathcal{L}}(f) \quad \text{where} \quad \widehat{\mathcal{L}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

Optimization Challenge

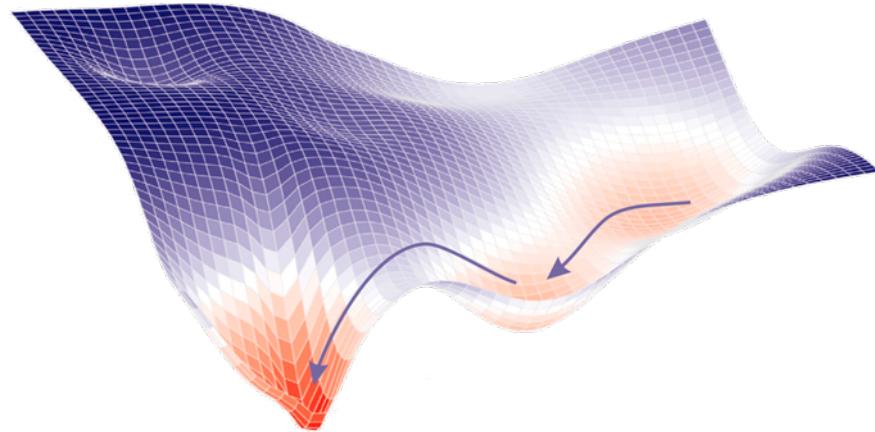
Our aim is typically finding an accurate model i.e.

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \widehat{\mathcal{L}}(f) \quad \text{where} \quad \widehat{\mathcal{L}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

- Challenge 2: How do you solve $\min_{f \in \mathcal{F}} \widehat{\mathcal{L}}(f)$?

Optimization

Goal: $\min_{f \in \mathcal{F}} \mathcal{L}(f)$



Some problems admit closed form solution (e.g. linear regression)



Modern ML problems are highly non-convex and non-linear.

Solution: Gradient descent

1. Workhorse of modern optimization
2. Optimize iteratively: Starting from an initial model f_0

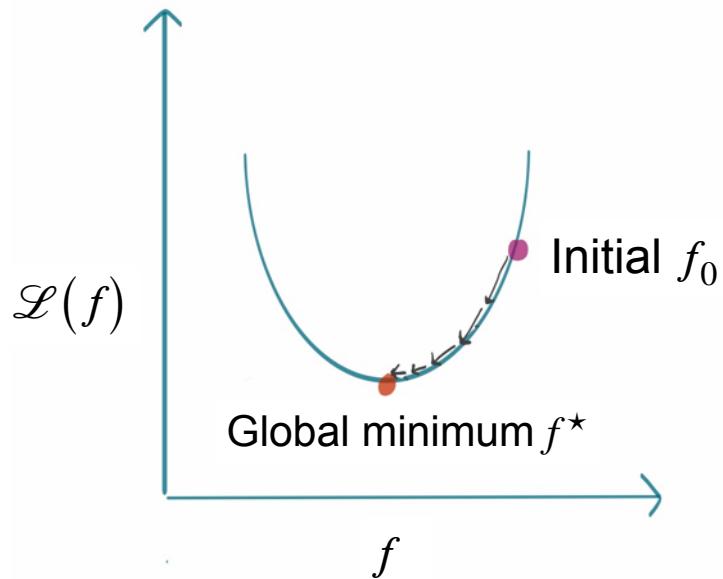
$$f_{i+1} = f_i - \eta \nabla \mathcal{L}(f_i)$$

Optimization

- Gradient descent iterations

$$f_{i+1} = f_i - \eta \nabla \mathcal{L}(f_i)$$

↑ ↑ ↗ ↗
New Current Learning Gradient
model model rate

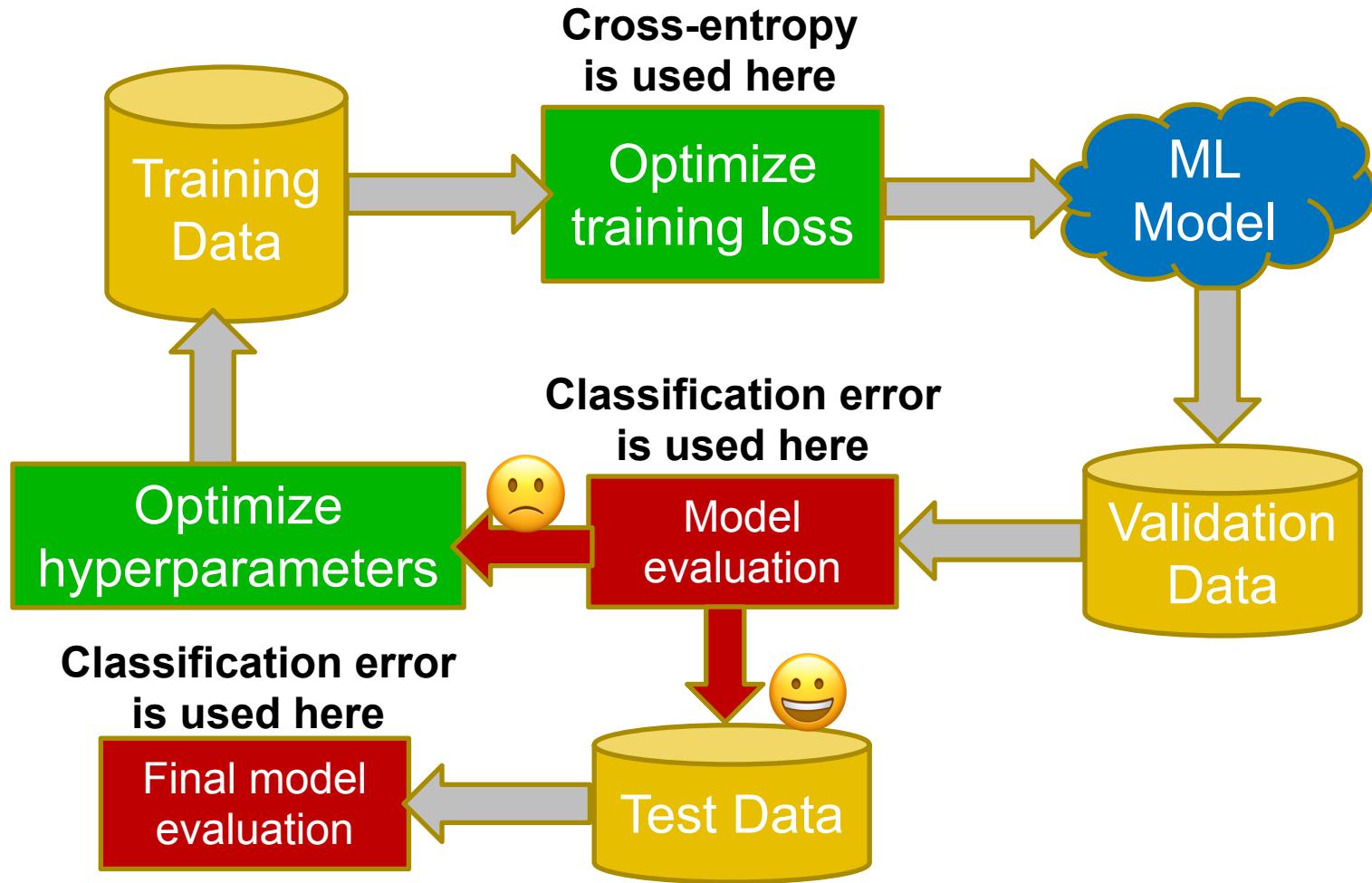


Back to loss functions

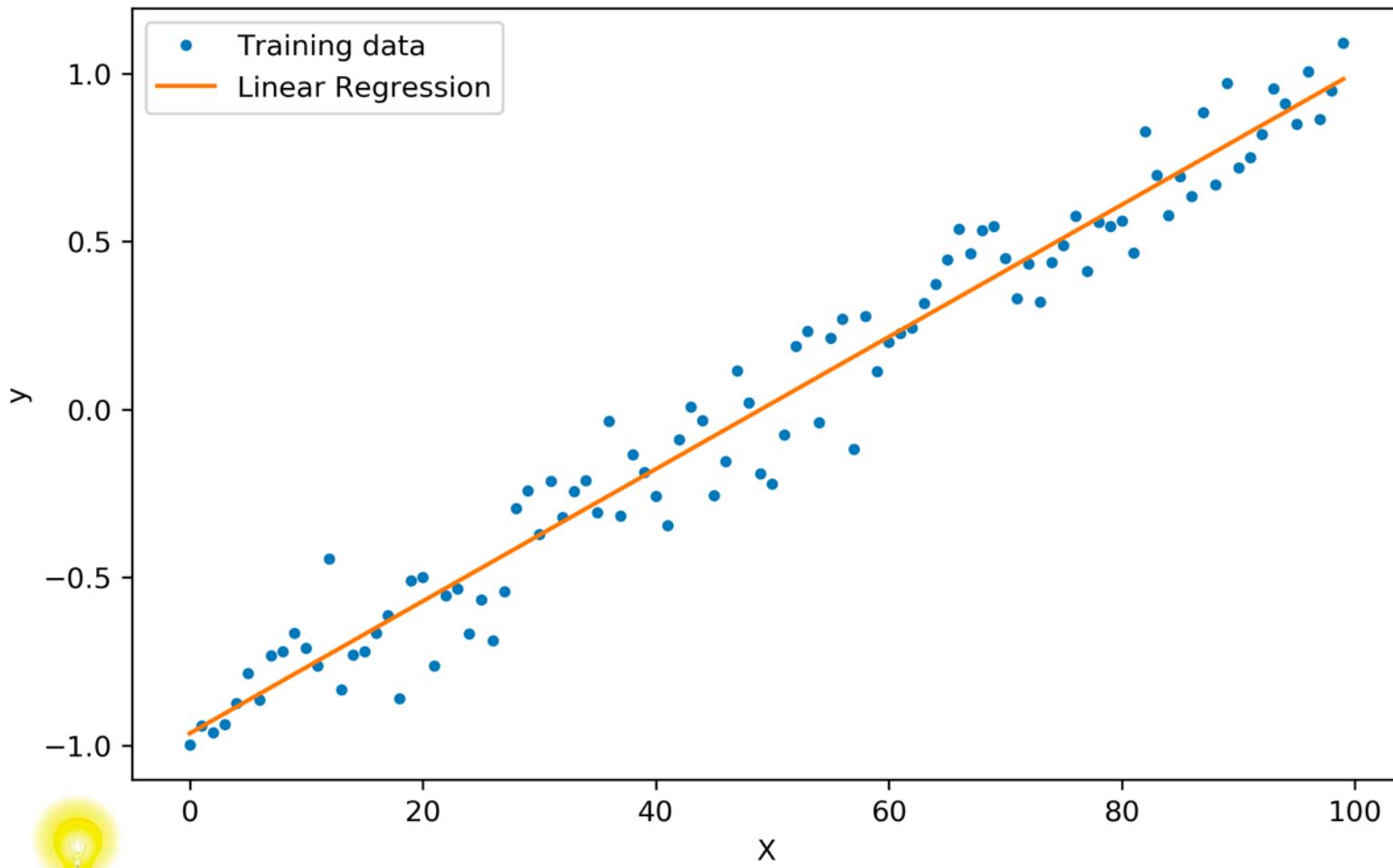
- Gradient descent requires gradient calculation
- Classification error $\mathbb{I}(\text{sign}(f(x)) \neq y)$ is not differentiable.
- Idea: Use a **surrogate loss function** instead
 1. Cross-entropy (logistic loss)
 2. Squared-loss
 3. ...
- In short: Training performance metric (training loss) may differ from test performance metric (test error).

Break

ML Pipeline



Optimization Fundamentals



Let us understand linear regression before deep learning.

Linear regression

1. Data: $(x_i, y_i)_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$
2. Model: Linear model $f(x; w) = x^\top w$
3. Squared Loss: $\ell(f(x), y) = \frac{1}{2} (f(x) - y)^2$
4. Optimize

$$\widehat{w} = \arg \min_{w \in \mathbb{R}^d} \widehat{\mathcal{L}}(w) \quad \text{where} \quad \widehat{\mathcal{L}}(w) = \frac{1}{2n} \sum_{i=1}^n (x_i^\top w - y_i)^2$$

Linear regression

- Let us move to matrix notation

$$X = \begin{matrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{matrix} \quad n \times d \text{ matrix} \qquad y = \begin{matrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{matrix} \quad n \times 1 \text{ vector}$$

- The loss function becomes

$$\widehat{\mathcal{L}}(w) = \frac{1}{2n} \sum_{i=1}^n (x_i^\top w - y_i)^2 = \frac{1}{2n} \|y - Xw\|_2^2$$

Linear regression

Optimization becomes

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{2} \|y - Xw\|_2^2$$

Solution can be found by differentiation.

Convex analysis

- Set the gradient to zero: $\nabla \mathcal{L}(w) = 0$

$$\nabla \mathcal{L}(w) = X^\top Xw - X^\top y$$

- Solution w^* satisfies: $X^\top Xw^* = X^\top y$

The normal equation

Linear regression

Solution satisfies: $X^\top X w^* = X^\top y$

- If $X^\top X$ is invertible (which requires $n \geq d$),
 - Unique solution $w^* = (X^\top X)^{-1} X^\top y$
 - **Intuition:** More equations than unknowns!
- Otherwise, **infinitely** many solutions.
 - Which solution to pick?
 - More on this later

Pseudo-inverse

A general solution: $w^\star = X^\dagger y$ where X^\dagger is pseudo-inverse

- $X^\dagger = (X^T X)^{-1} X^T$ if columns of X are linearly independent

$$n \geq d$$

Over-determined
problem

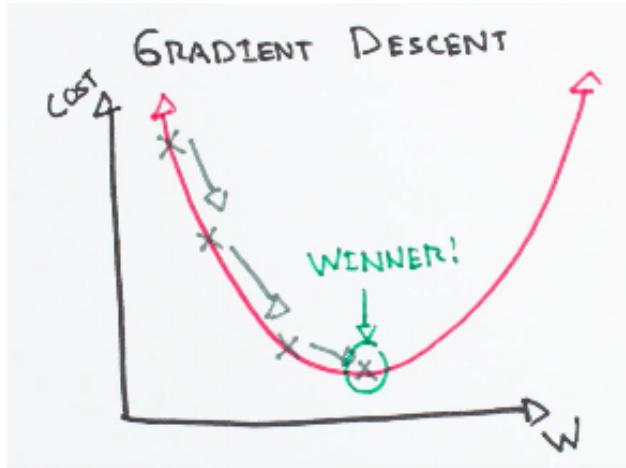
- $X^\dagger = X^T (X X^T)^{-1}$ if rows are linearly independent

$$n < d$$

Over-parameterized
problem

Linear regression

- Pseudo-inverse is simple and closed-form
- But matrix inversion is an expensive operation
 - $w^* = (X^T X)^{-1} X^T y$ requires $O(d^3 + dn)$ operations naively
Matrix inversion: $O(d^3)$ Matrix-vector multiplication: $O(nd)$
- How to solve linear regression without inversion?



Linear regression

- Recall the gradient:

$$\nabla \mathcal{L}(w) = X^\top X w - X^\top y$$

- Fix some initial weights w_0
- Fix a learning rate $\eta > 0$
- Gradient iterations are given by

$$w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t)$$

Linear regression

- Gradient iterations are given by

$$w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t)$$

- Plugging in the gradient formula

$$w_{t+1} = w_t - \eta (X^\top X w_t - X^\top y)$$

- Computational cost of each iteration

$$O(nd)$$

Only need matrix-vector multiplication

Linear regression

Plugging in the gradient formula

$$w_{t+1} = w_t - \eta(X^\top X w_t - X^\top y)$$

Questions:

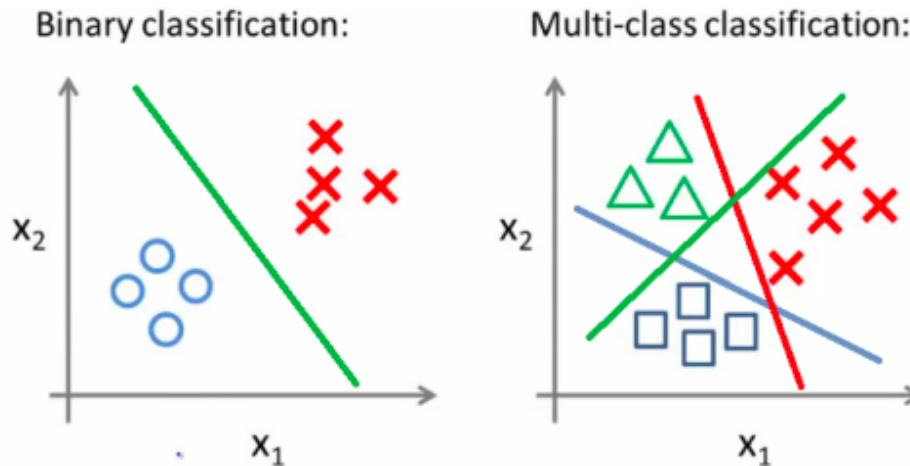
- How many iterations needed?
- Which learning rate?
- Can we speed up further?

To be discussed in
future lectures

Answers are important for **deep learning** as well.

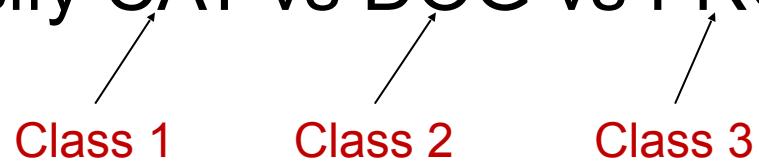
Multiclass Classification

- So far, we assumed label is a scalar
- Modern ML models can classify 100s of objects



- How do we handle multiple classes?

One-hot encoding

- For a K -class problem, use K dimensional labels
 - Ex: Classify CAT vs DOG vs FROG.
 - CAT receives label [1 0 0]
 - DOG receives label [0 1 0]
 - FROG receives label [0 0 1]
- 
- Class 1 Class 2 Class 3

Example Dataset

Inputs x_i



Labels y_i

[0 0 1]



[0 1 0]



[1 0 0]



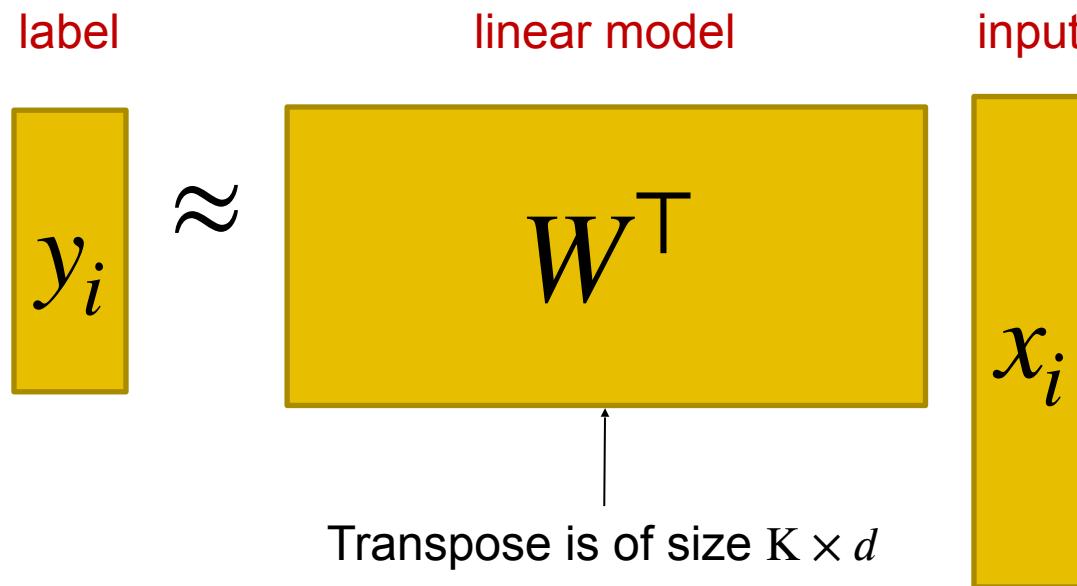
[0 1 0]

Multiclass Classifier

- Suppose $x \in \mathbb{R}^d$, $y \in \mathbb{R}^K$
- Need a ML model $f: \mathbb{R}^d \rightarrow \mathbb{R}^K$
- What does this mean for linear model?

Multiclass Linear Model

- Multiclass linear model
 - a **matrix** W rather than **vector** w
 - dimensions $d \times K$



Multiclass Optimization

- Back to linear regression with squared-loss

$$\widehat{\mathcal{L}}(W) = \frac{1}{2} \sum_{i=1}^n \|y_i - W^\top x_i\|_2^2$$

The L_2 distance between
label vector and prediction

Multiclass Optimization

- Back to linear regression with squared-loss

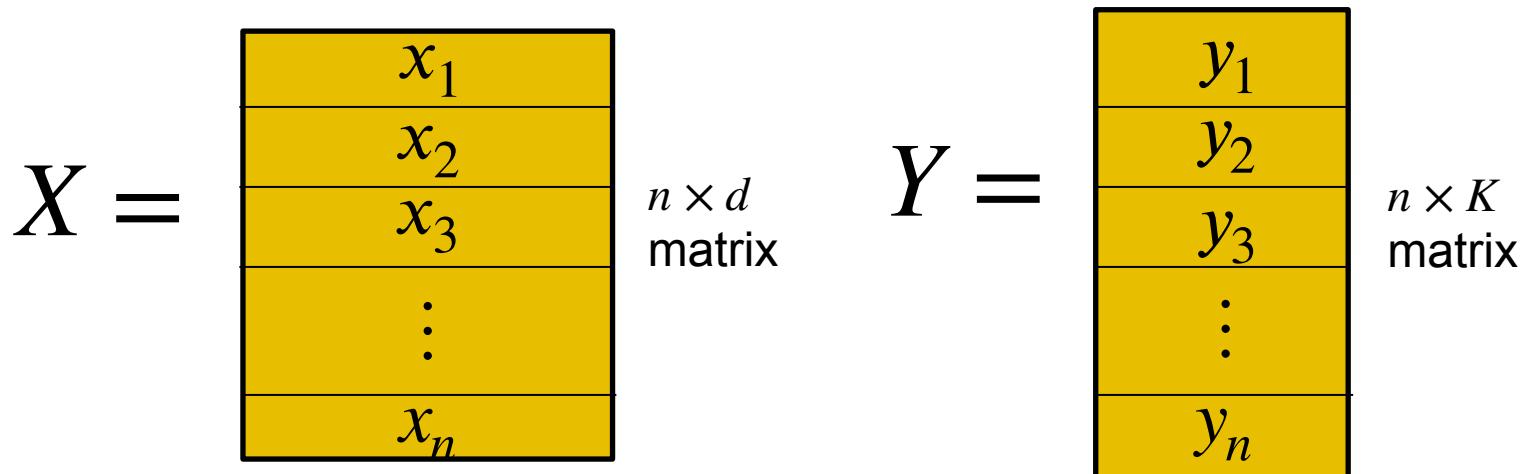
$$\begin{aligned}\widehat{\mathcal{L}}(W) &= \frac{1}{2} \sum_{i=1}^n \|y_i - W^\top x_i\|_2^2 \\ &= \frac{1}{2} \|Y - XW\|_F^2\end{aligned}$$

$$\boxed{\|Z\|_F = \sqrt{\sum_{i,j} Z_{i,j}^2}}$$

Multiclass Optimization

- Back to linear regression with squared-loss

$$\widehat{\mathcal{L}}(W) = \frac{1}{2} \sum_{i=1}^n \|y_i - W^\top x_i\|_2^2 = \frac{1}{2} \|Y - XW\|_F^2$$



Optimizing with Gradient Descent

- Gradient is given by

$$\nabla \mathcal{L}(W) = X^T(XW - Y)$$

- Plugging in the gradient formula

$$W_{t+1} = W_t - \eta(X^T X W_t - X^T Y)$$

- Computational cost of each iteration

$$O(K \times Nd)$$

- The equations remain the same! 😊
- To be discussed: **cross-entropy loss**

Summary

1. Machine learning fundamentals
2. Gradient descent basics
 - Linear regression
3. Multiclass models & one-hot encoding
4. Next class: Getting started with neural nets