[544] PyTorch Optimization

Meenakshi Syamkumar

Learning Objectives

- write a PyTorch optimization loop to find inputs that minimize/maximize an output
- frame model training as an optimization problem minimizing loss
- prepare datasets using DataSet and DataLoader from sources like CSVs

Outline

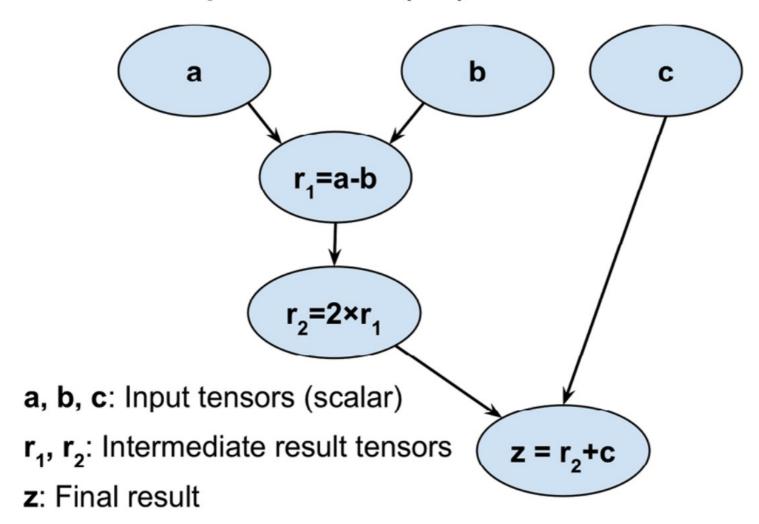
Optimization

- Calculations as DAGs
- Iterative approach

Machine Learning

- Brief background
- Machine Learning as Optimization

Computation graph implementing the equation $z = 2 \times (a-b) + c$



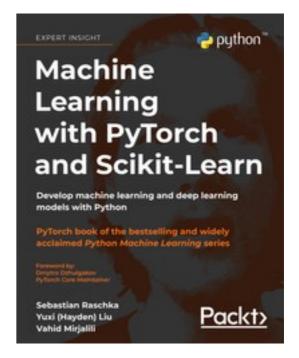


Figure 13.1: How a computation graph works

PyTorch can calculate how small changes in one variable in the DAG impacts another. Example: if b *increases* by 0.001, z will *decrease* by 0.002. The *gradient* of z with respect to b is -2.

Optimization: if we want z to be large, decreasing b a little (how much?) is probably a good idea.

Making a small improvement

```
a = torch.tensor(3.0)
b = torch.tensor(4.0)
c = torch.tensor(5.0)
z = 2 * (a - b) + c
```

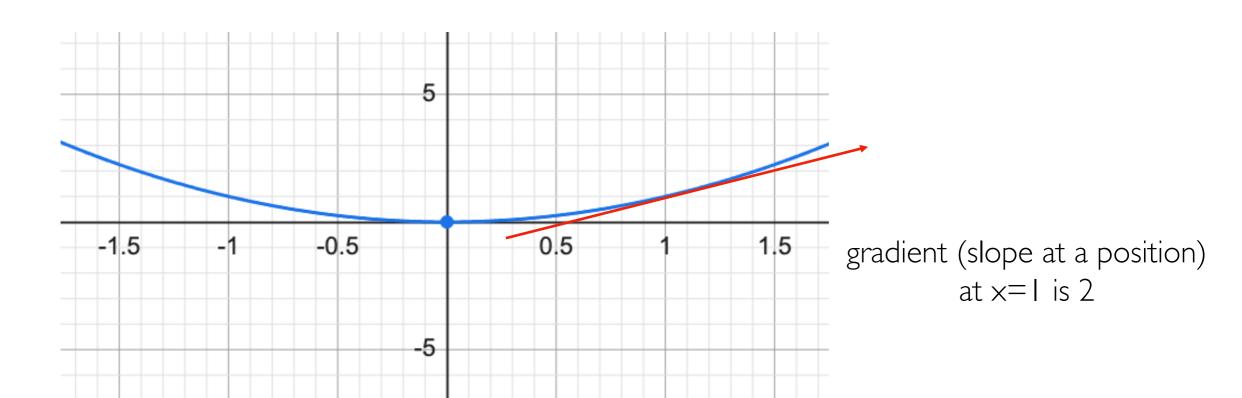
Scenario

- We want z to be large
- We're allowed to change b (but to what?)

Tracking gradients

```
a = torch.tensor(3.0)
b = torch.tensor(4.0, requires_grad=True)
c = torch.tensor(5.0)

z = 2 * (a - b) + c
```



Calculating gradients

Accumulating gradients

```
a = torch.tensor(3.0)
b = torch.tensor(4.0, requires_grad=True)
c = torch.tensor(5.0)

z = 2 * (a - b) + c
z.backward()
b.grad

z = 2 * (a - b) + c
z.backward()
b.grad

z = 2 * (a - b) + c
z.backward()
b.grad
```

careful, gradients accumulate in .grad everytime you call backward (has uses, but not usually what we want)

Taking steps

```
a = torch.tensor(3.0)
b = torch.tensor(4.0, requires_grad=True)
c = torch.tensor(5.0)

optimizer = ????

z = 2 * (a - b) + c
z.backward()
b.grad -2

optimizer.step()
```

step() will make b a little bigger or a little smaller, depending on gradient, and whether we're minimizing or maximizing

Stochastic Gradient Descent (SGD) Optimizer

```
a = torch.tensor(3.0)
                                            what I to be big
b = torch.tensor(4.0, requires grad=True)
c = torch.tensor(5.0)
optimizer = torch.optim.SGD([b], maximize=True,
                                lr=0.1)
                                  learning rate specifies how
                                 much step() should change b
z = 2 * (a - b) + c
z.backward()
b.grad
                    \rightarrow b += b.grad * lr
optimizer.step()
                                -2 * 0.1 = -0.2
  now b is 3.8
                          (use -= if minimizing)
```

Clearing gradients (to prep for another step)

```
a = torch.tensor(3.0)
b = torch.tensor(4.0, requires grad=True)
c = torch.tensor(5.0)
optimizer = torch.optim.SGD([b], maximize=True,
                             lr = 0.1)
z = 2 * (a - b) + c
z.backward()
b.grad
optimizer.step()
optimizer.zero grad()
b.grad
```

Iteratively improving

```
a = torch.tensor(3.0)
b = torch.tensor(4.0, requires grad=True)
c = torch.tensor(5.0)
optimizer = torch.optim.SGD([b], maximize=True,
                                lr = 0.1)
                                each iteration of optimization
for epoch in range (10):
                                   is called an "epoch"
     z = 2 * (a - b) + c
     z.backward()
     optimizer.step()
     optimizer.zero grad()
              many small improvements have been made
print(b)
```

Demos...

Outline

Optimization

- Calculations as DAGs
- Iterative approach

Machine Learning

- Brief background
- Machine Learning as Optimization

Machine Learning, Major Ideas

Categories of Machine Learning:

- Reinforcement learning: agent makes series of actions to maximize reward
- Unsupervised learning: looking for general patterns
- Supervised learning: train models to predict unknowns (today)

Models are functions that return predictions:

Example:

```
def weather_forecast(temp_today, temp_yesterday):
    ...
    return temp_tomorrow
```

Machine Learning, Major Ideas

Categories of Machine Learning:

- Reinforcement learning: agent makes series of actions to maximize reward
- Unsupervised learning: looking for general patterns
- Supervised learning: train models to predict unknowns (today)

Models are functions that return predictions:

```
def my_model(some_info):
    ...
    return some_prediction
```

computation usually involves some calculations (multiply, add) with various numbers (parameters). Training is finding parameters that result in good predictions for known training data

Example:

```
def weather_forecast(temp_today, temp_yesterday):
    ...
    return temp_tomorrow
```

Goal: Learning from Data

	x1	x2	У
0	2	8	5
1	9	2	6
2	4	1	0
3	7	9	7
4	2	2	3
5	3	4	3
6	3	5	9
7	7	1	4
8	6	6	3
9	4	3	?
10	1	2	?
11	2	9	?

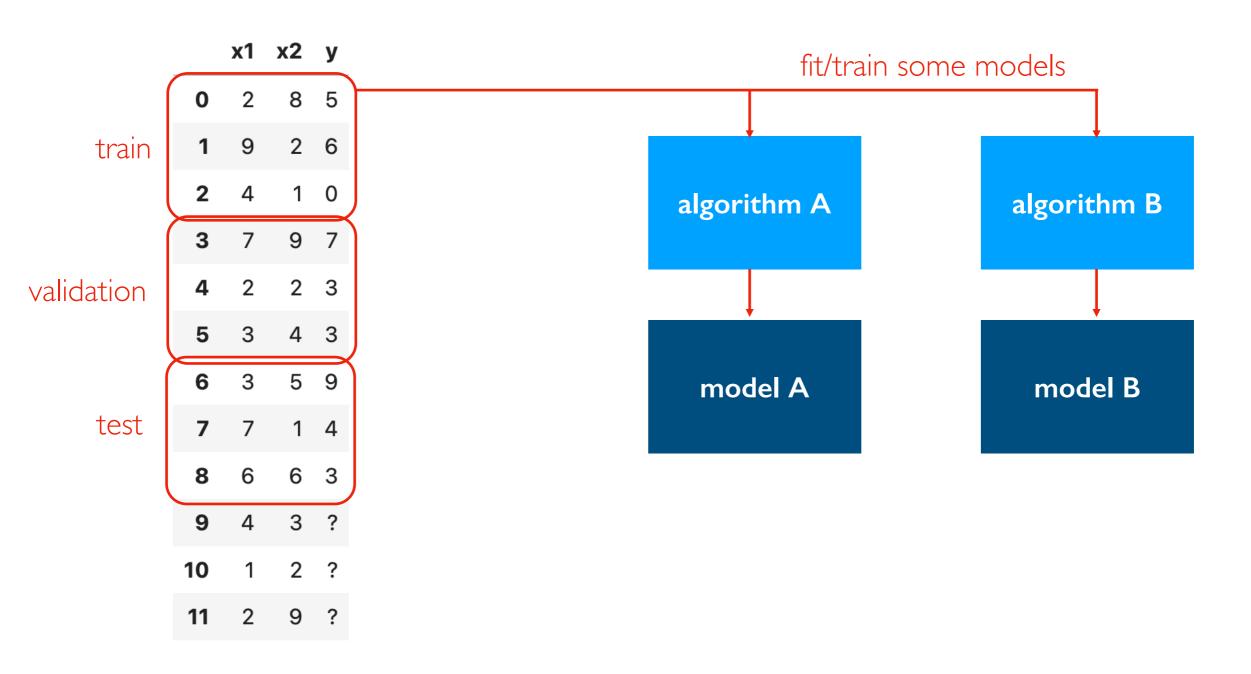
how can the cases where we DO know y help us predict the cases where we do not?

Split Known Cases

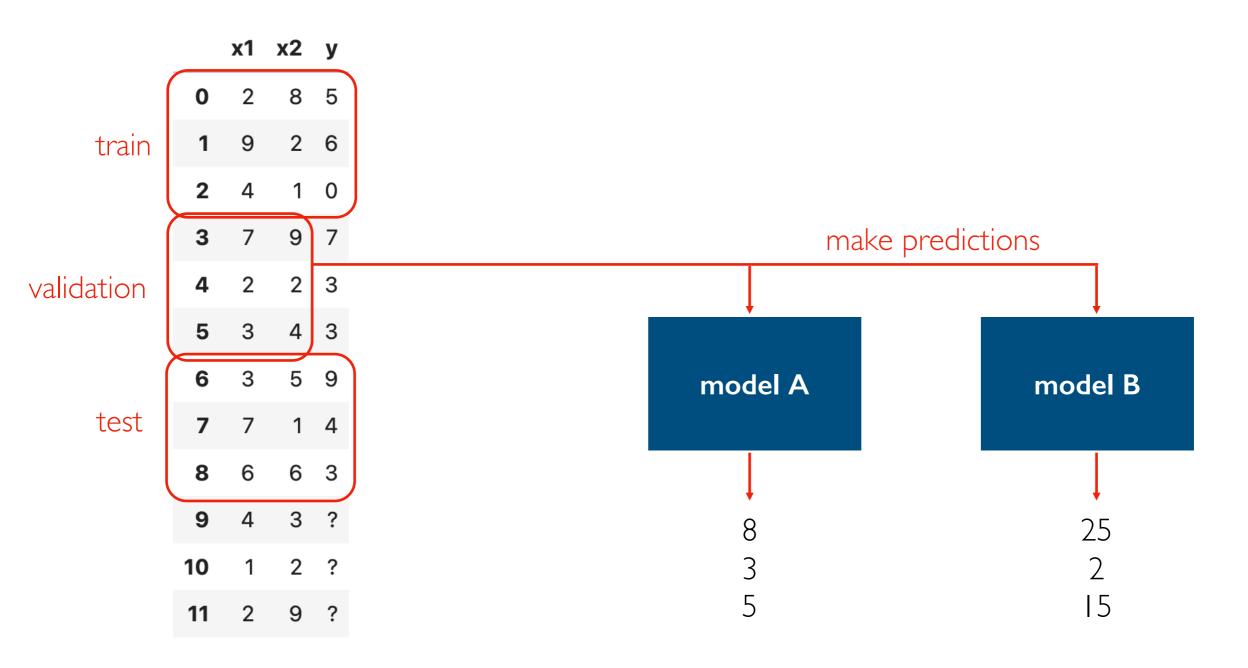
		х1	x2	у
	0	2	8	5
train	1	9	2	6
	2	4	1	0
	3	7	9	7
validation	4	2	2	3
	5	3	4	3
	6	3	5	9
test	7	7	1	4
	8	6	6	3
	9	4	3	?
	10	1	2	?
	11	2	9	?

random split

Train Models



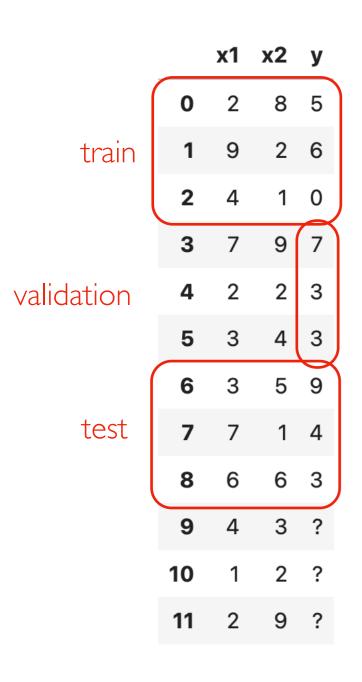
Predict with Models



Measure Loss

which model

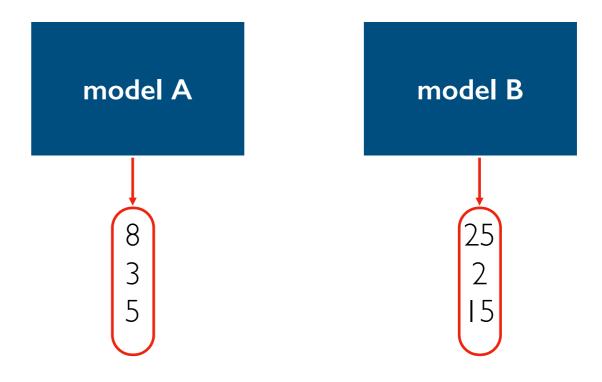
predicts better?



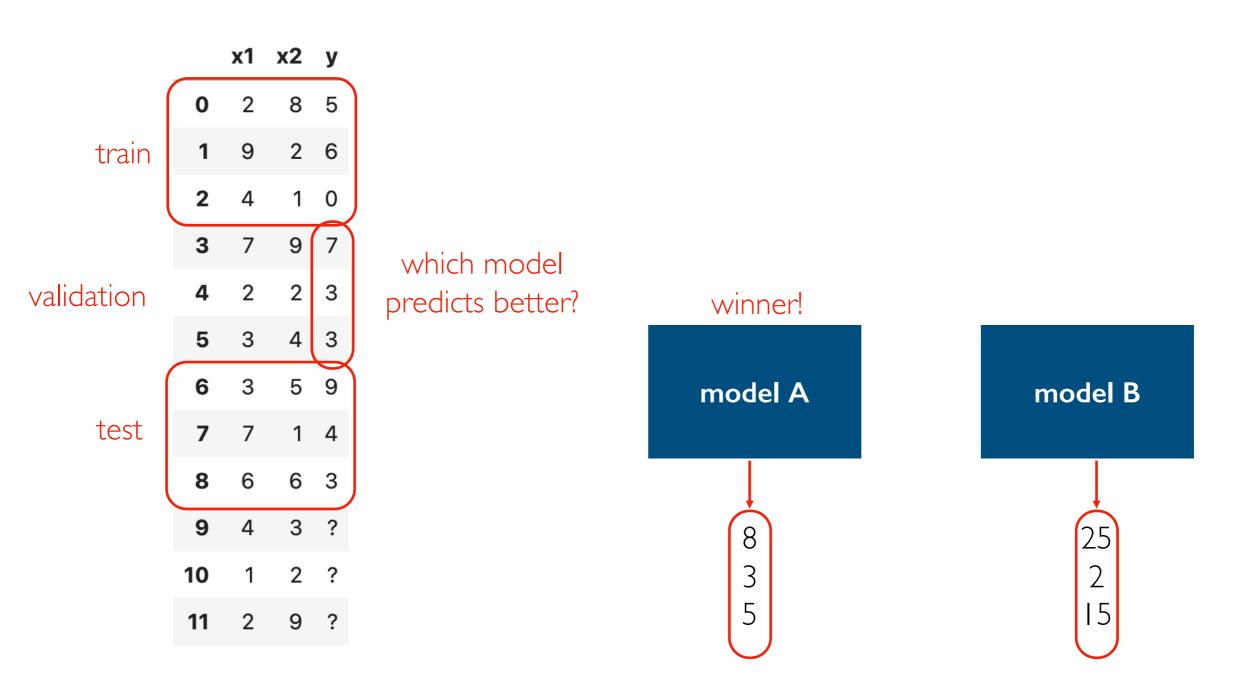
у	Р	err	err^2	
7	8	I		
3	3	0	0	
3	5	2	4	
•				

MSE (mean squared error) is 5/3 = 1.666

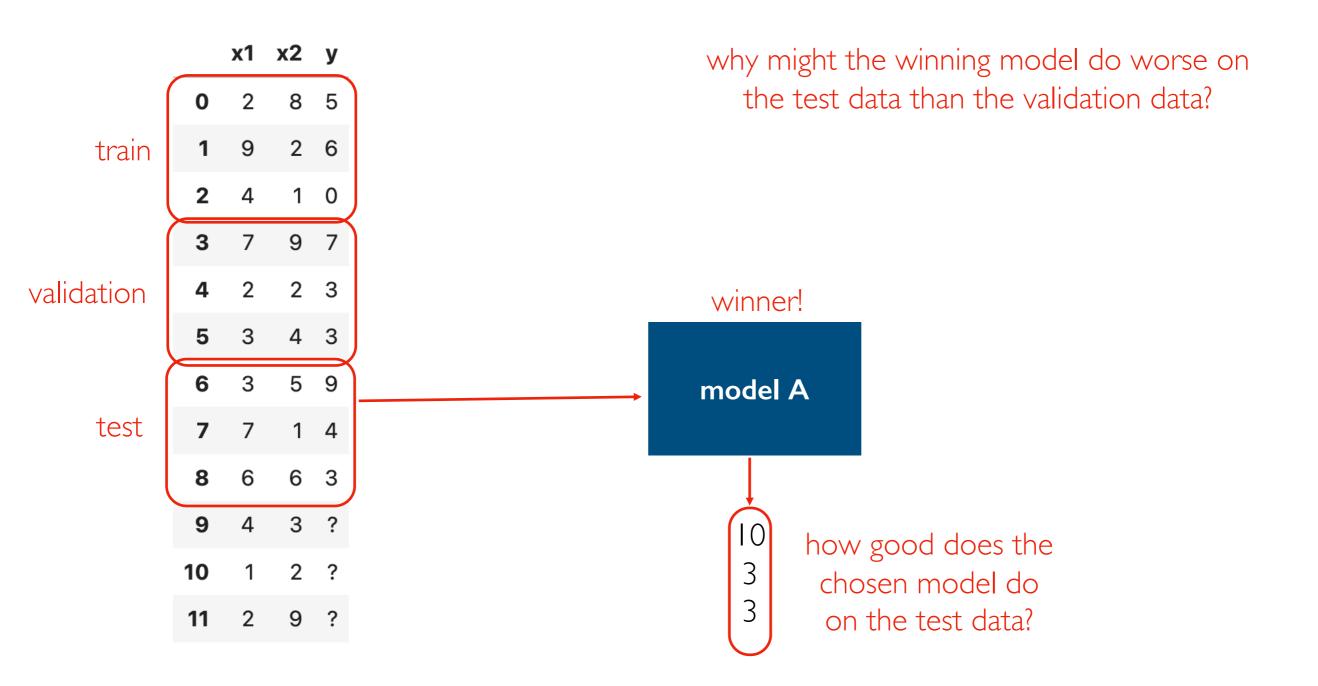
Loss functions measure how bad predictions are (MSE is one possible metric)



Choose best model

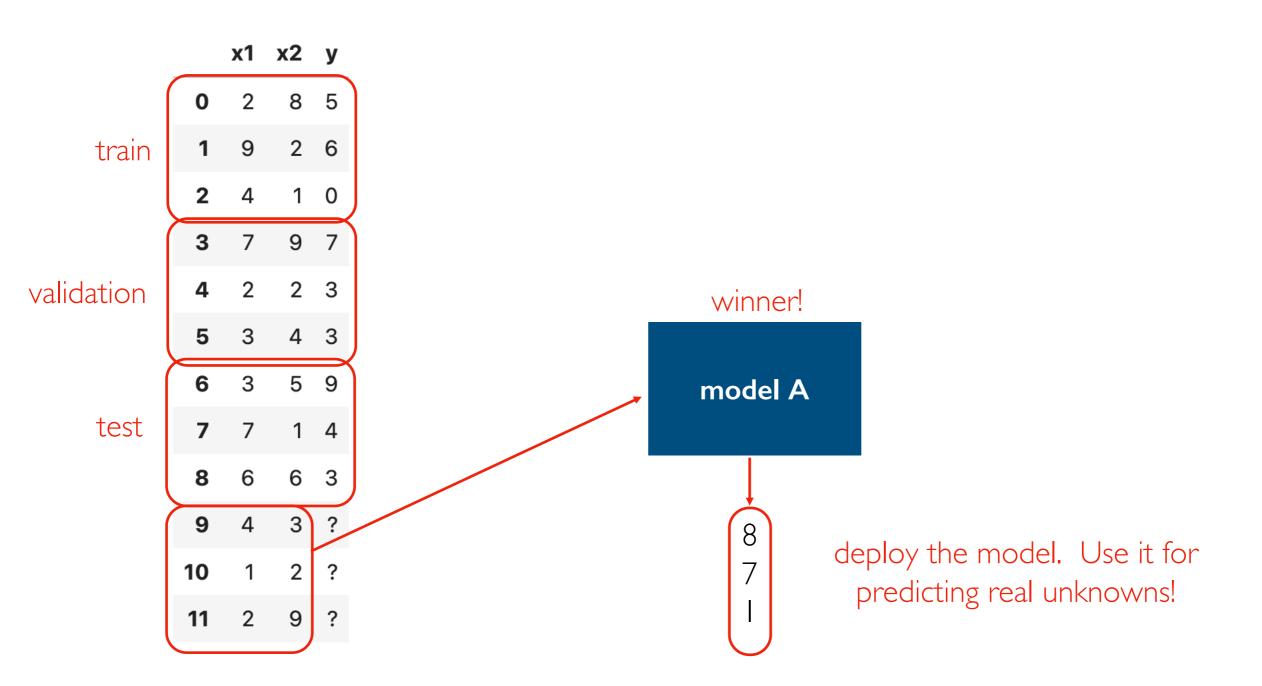


How does the winner do on something new?



models that do good on train data but bad on validation/test data have "overfitted"

Deploy!



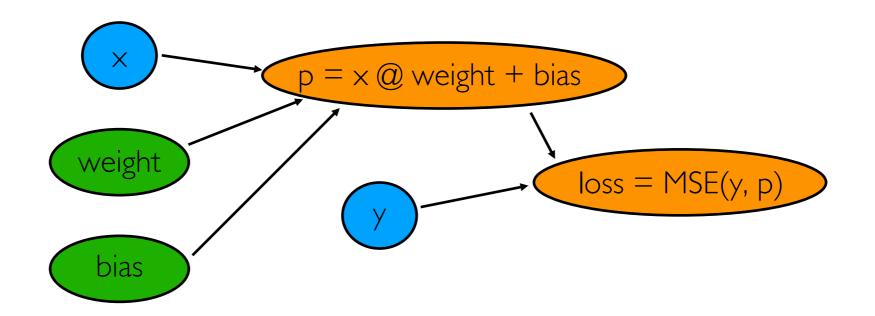
Outline

Optimization

- Calculations as DAGs
- Iterative approach

Machine Learning

- Brief background
- Machine Learning as Optimization

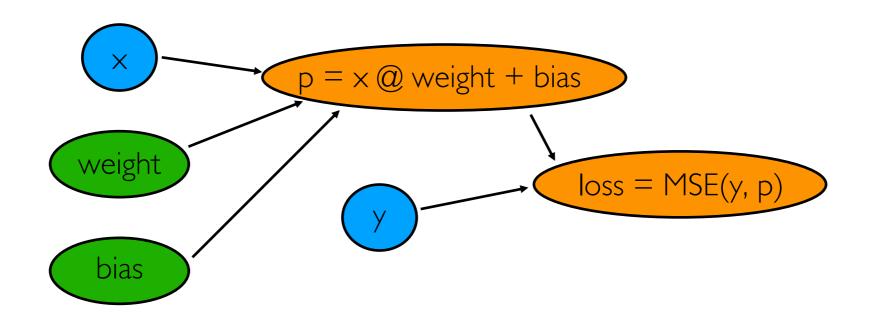


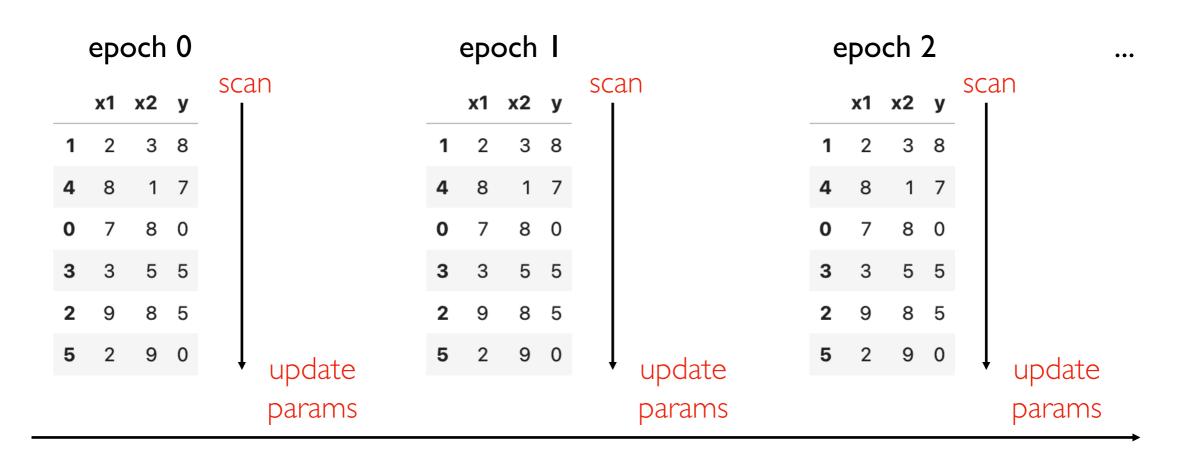
x and y are known (these are matrices/vectors).
what should weight and bias (parameters) be?

def model(data): return data @ weight + bias

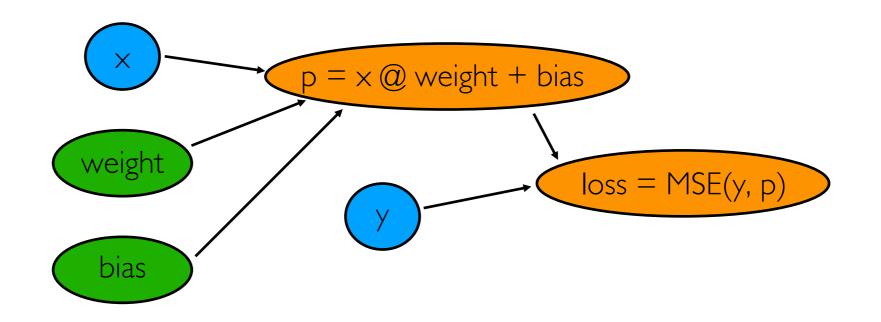
p = model(x)loss = MSE(y, p) MSE (means squared error) measures how different predictions are from real values, so we want small loss (optimization).

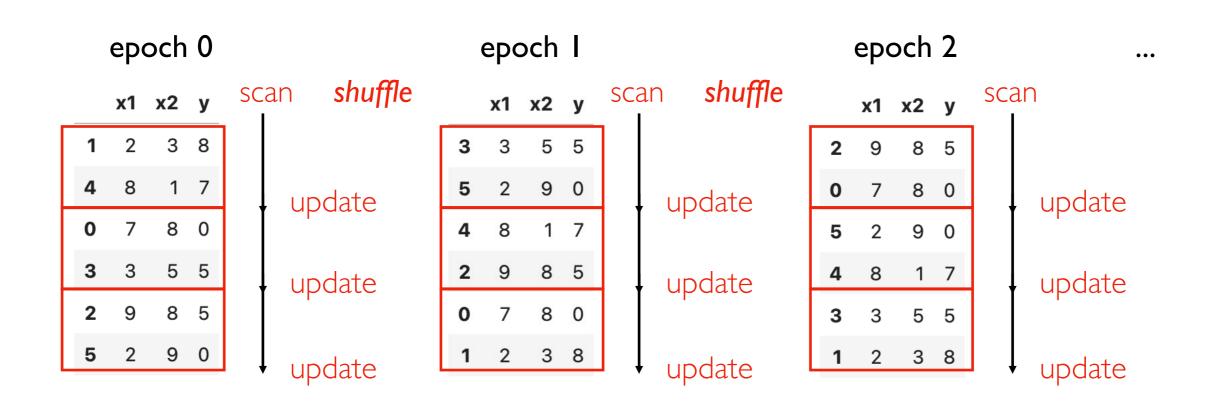
If gradient of loss with respect to weight is positive, then decrease weight.





gradient descent. slow (consider all data each update), and data might not fit in RAM





stochastic gradient descent. shuffle each time, process in small batches that fit in memory

Demos...