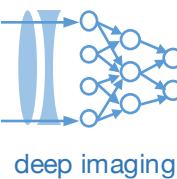


Lecture 6: Ingredients for Machine Learning

Machine Learning and Imaging

BME 590L
Roarke Horstmeyer



Outline

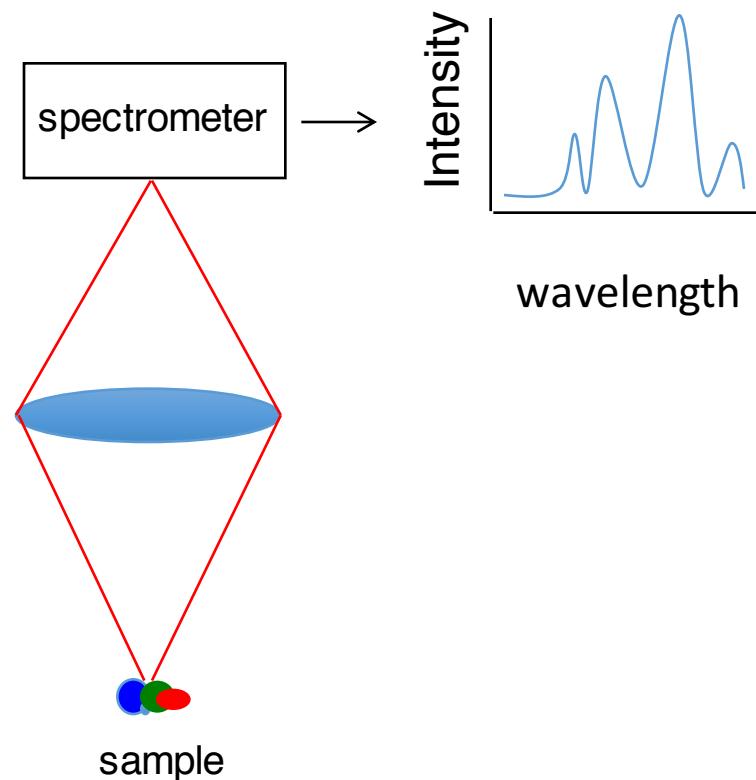
- Review spectral unmixing (last class)
- From optimization to machine learning
- Ingredients for ML
- Example: linear classification of images
 - Train/test data
 - Linear regression model
 - 3 ways to solve

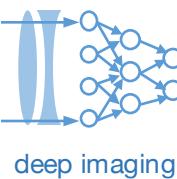
Last time: simple example of spectral unmixing

(For whatever reason, whenever I get confused about optimization, I think about this example)

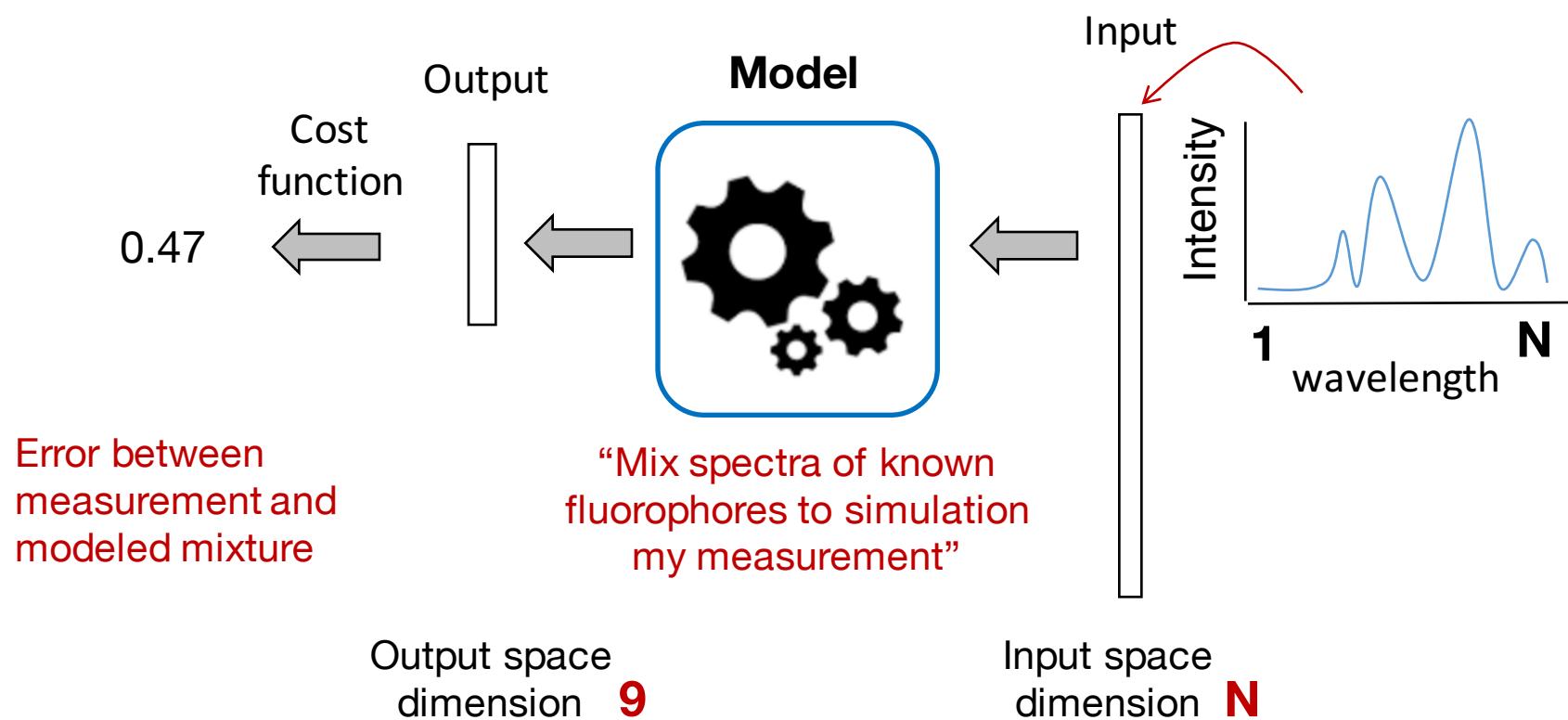
The setup:

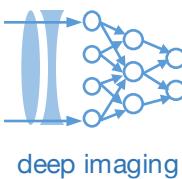
- measure the color (spectral) response of a sample (e.g., how much red, green and blue there is, or several hundred measurements of its different colors).
- You know that the sample can only contain 9 different fluorophores.
- What % of each fluorophores is in your sample?



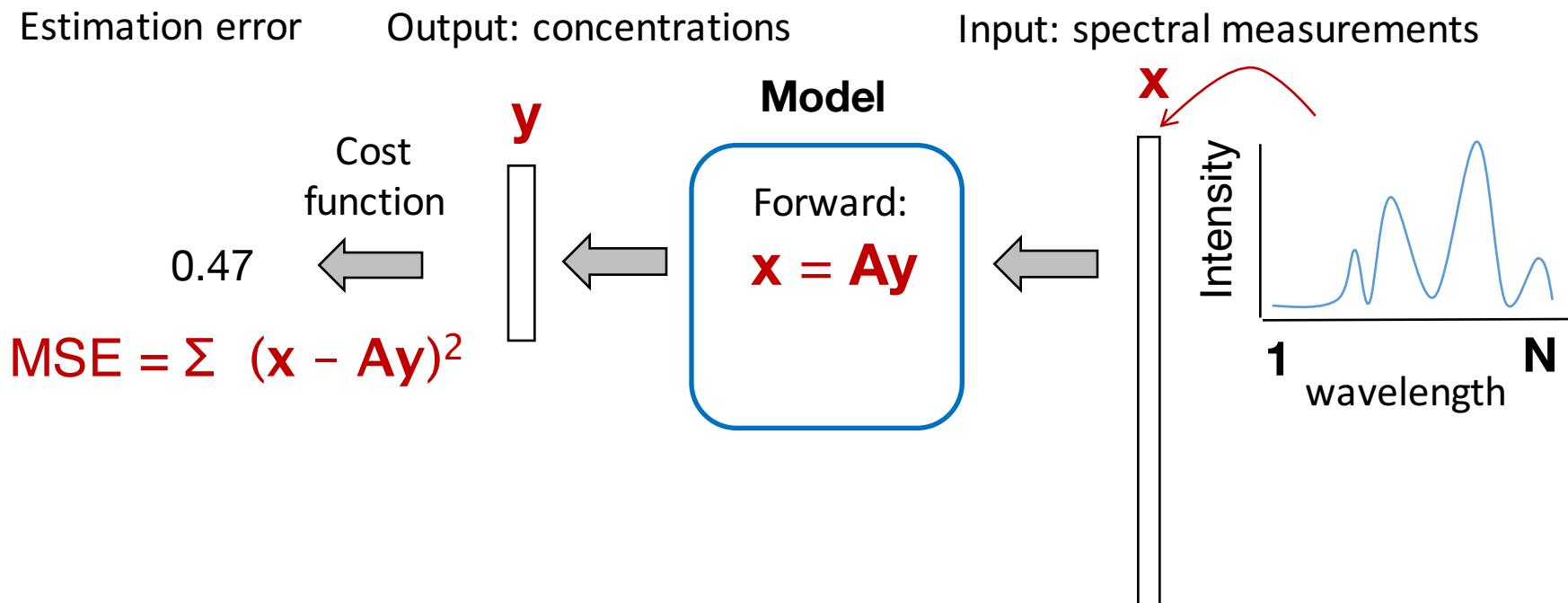


Optimization pipeline for spectral unmixing



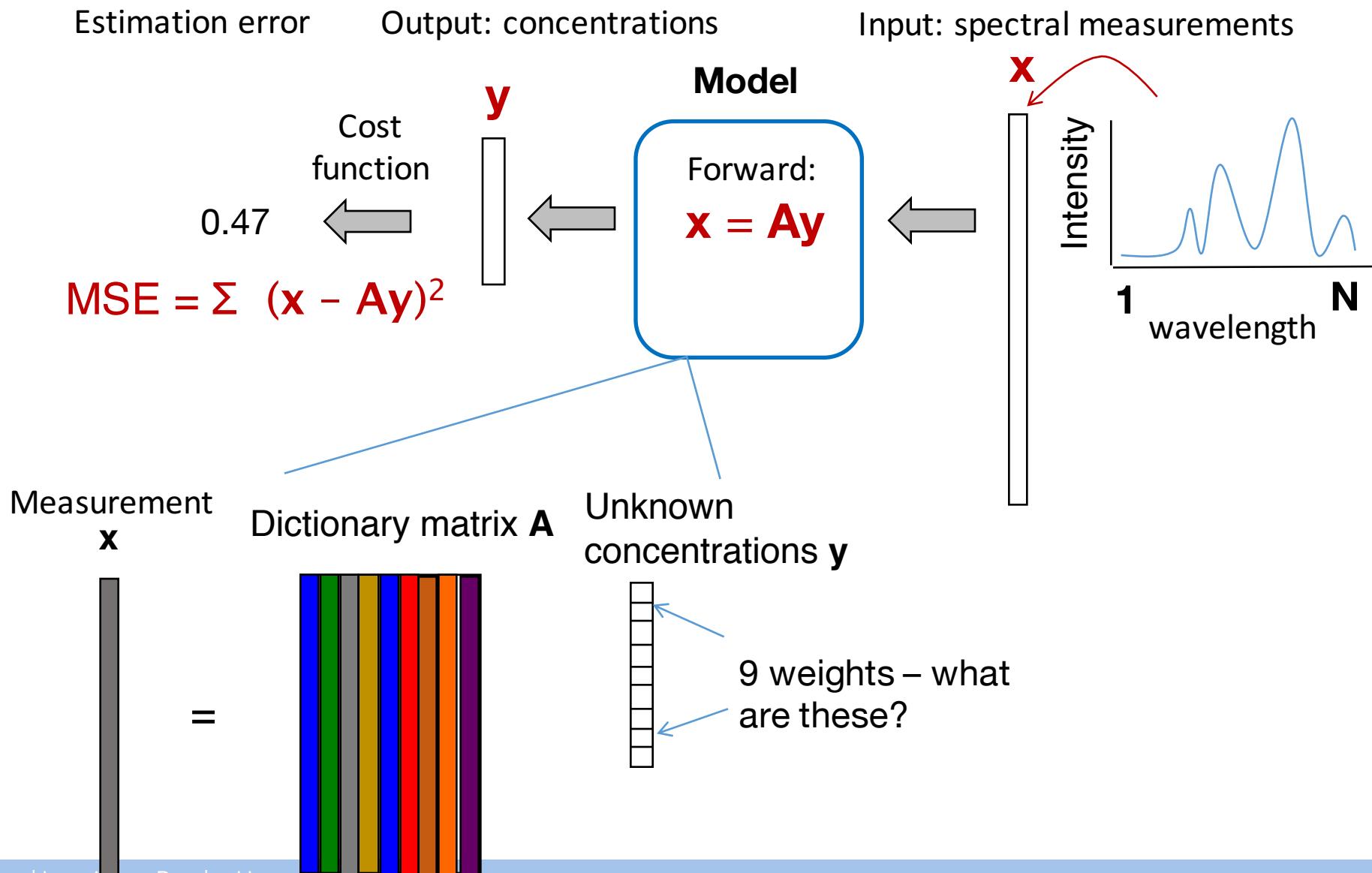


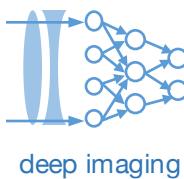
Optimization pipeline for spectral unmixing



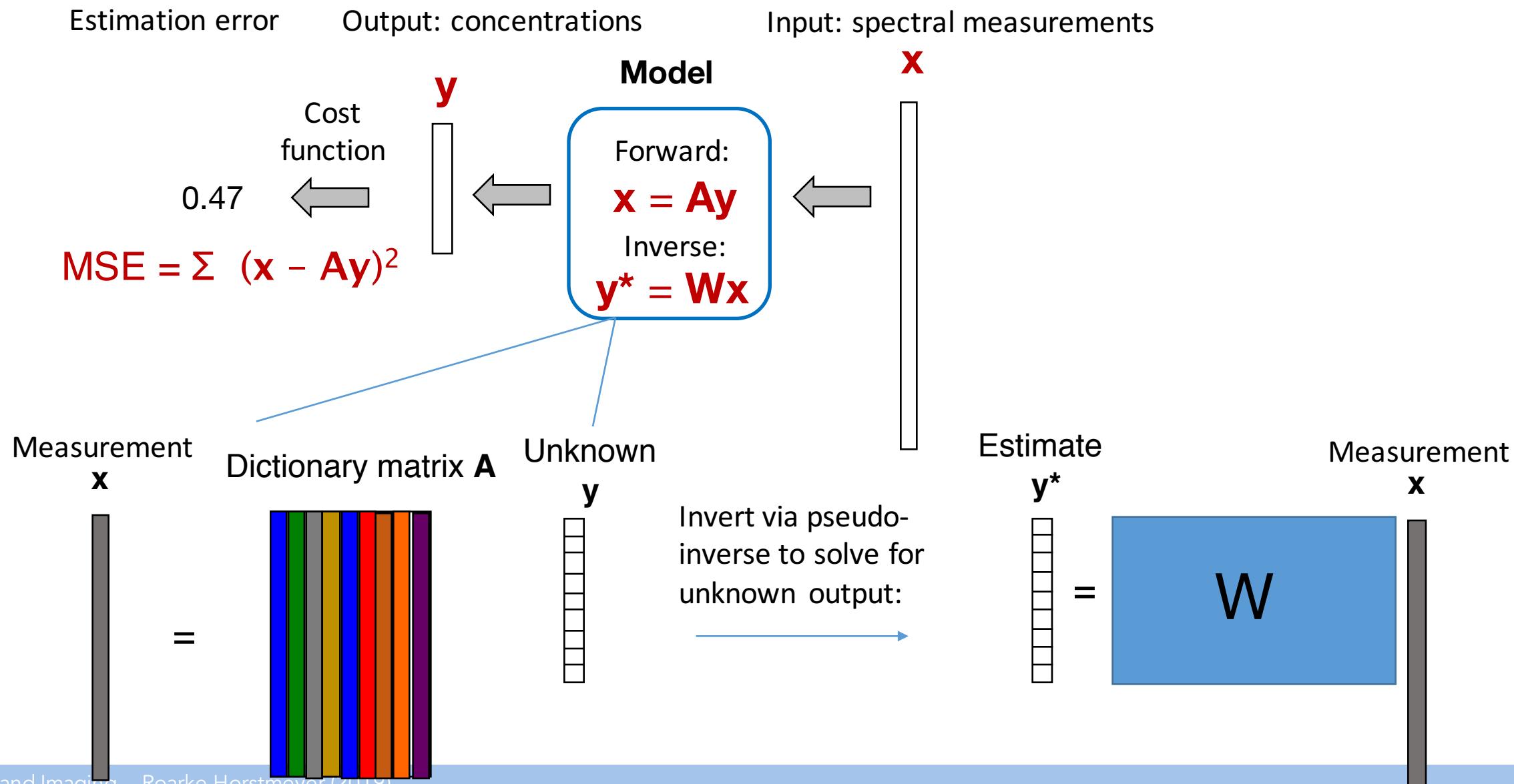
*Note: notation changed from last time to be consistent with we'll use in the future

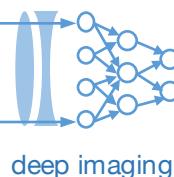
Optimization pipeline for spectral unmixing



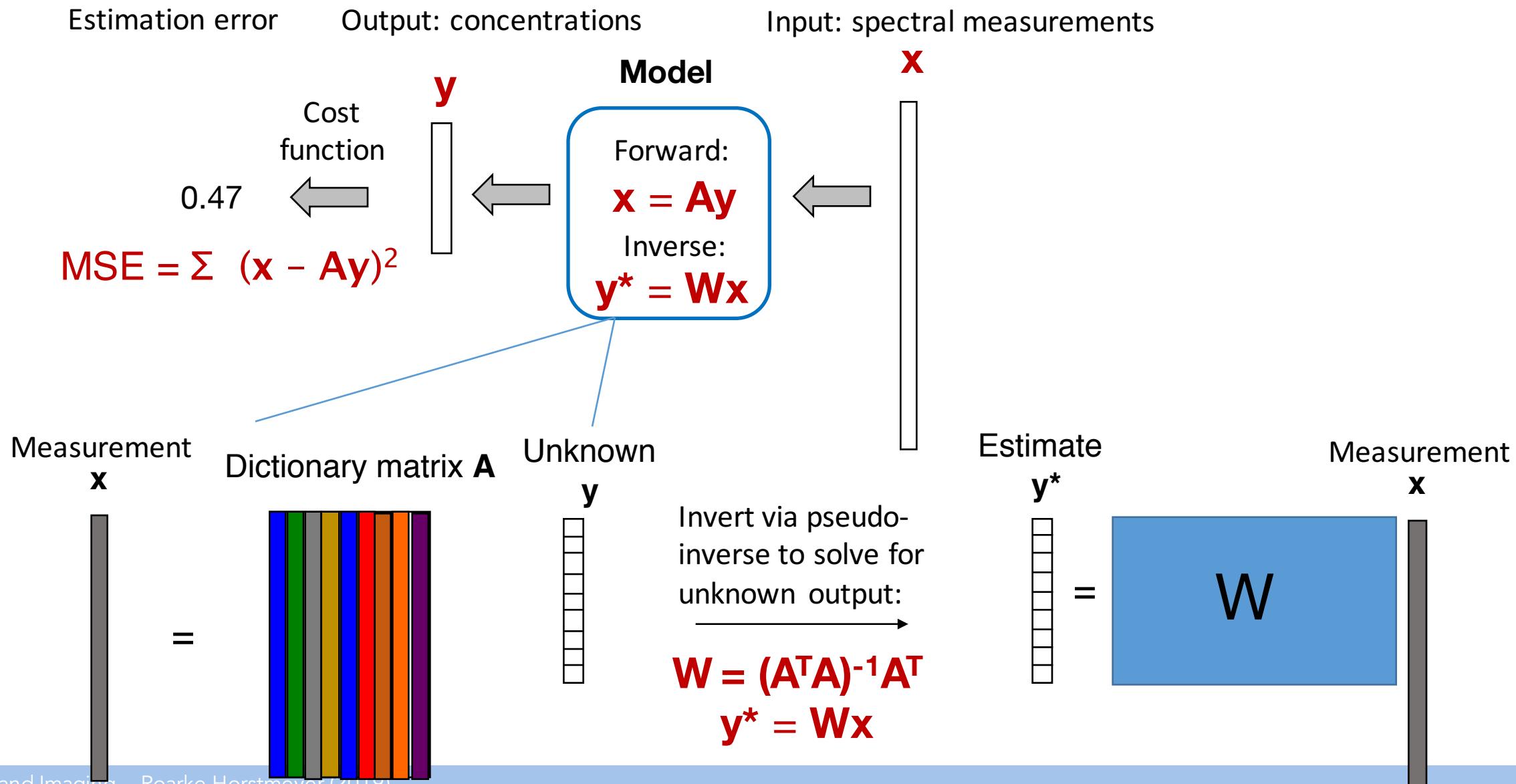


Optimization pipeline for spectral unmixing

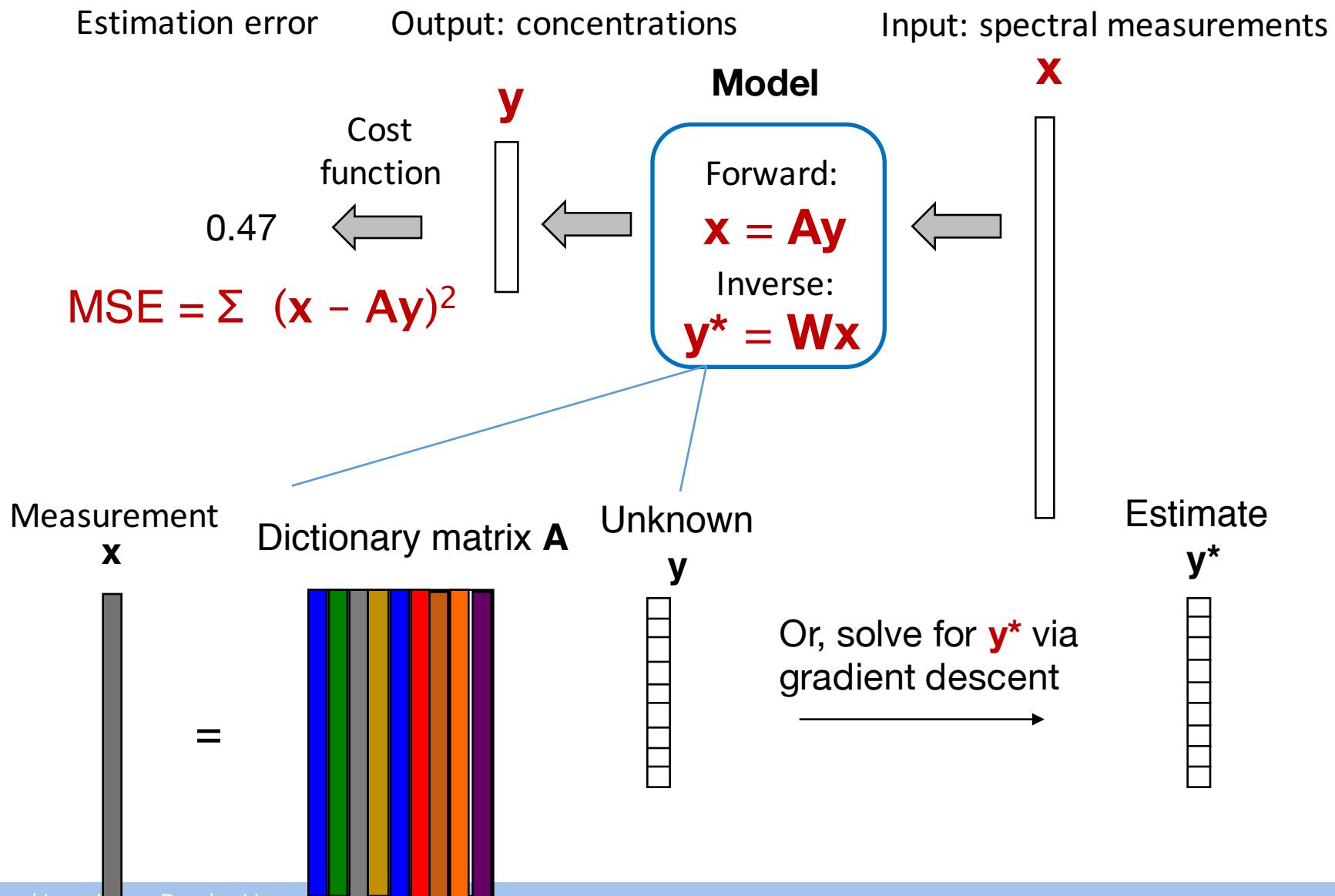




Optimization pipeline for spectral unmixing



Optimization pipeline for spectral unmixing

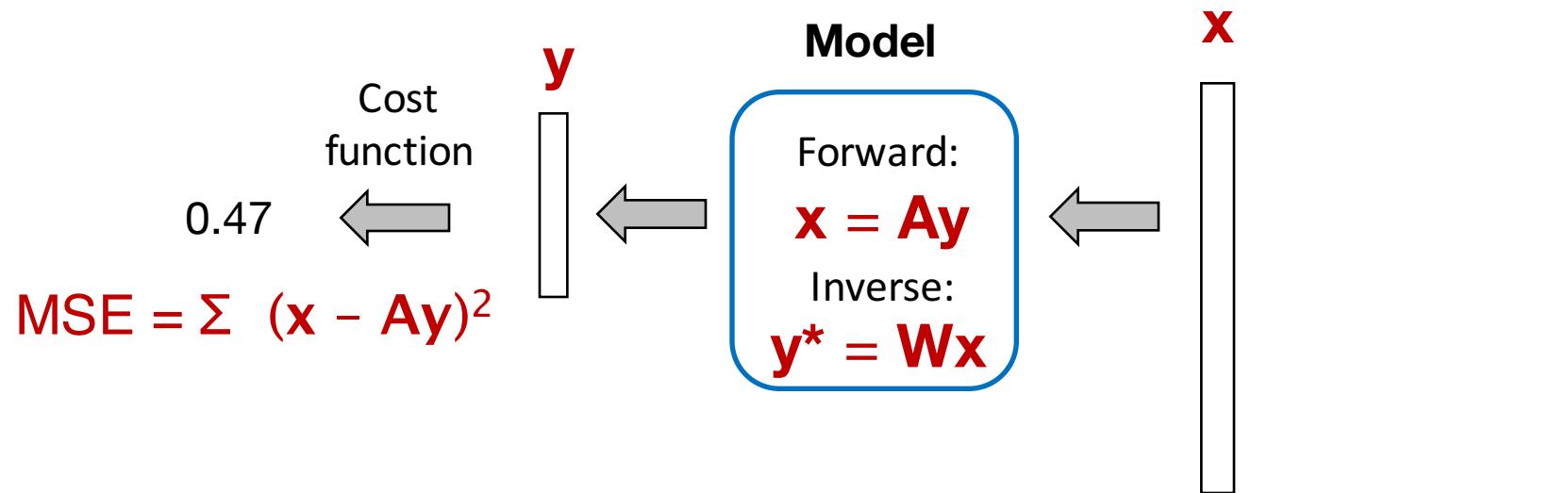


Optimization pipeline for spectral unmixing

Estimation error

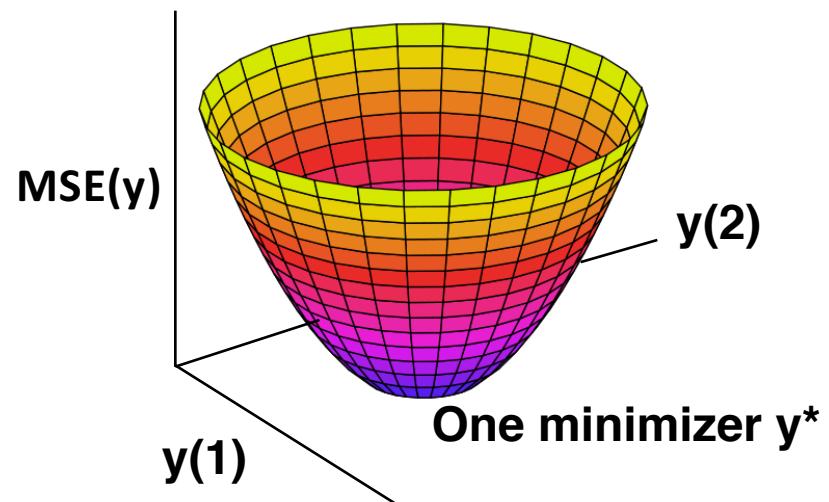
Output: concentrations

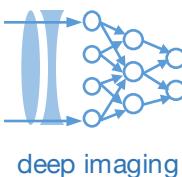
Input: spectral measurements



$$MSE(\mathbf{y}) = \sum \text{spectral measurements} (\mathbf{x} - \mathbf{A}\mathbf{y})^2$$

$$\frac{d}{dy} MSE = \mathbf{A}^T(\mathbf{x} - \mathbf{A}\mathbf{y})$$





Optimization pipeline for spectral unmixing

Estimation error

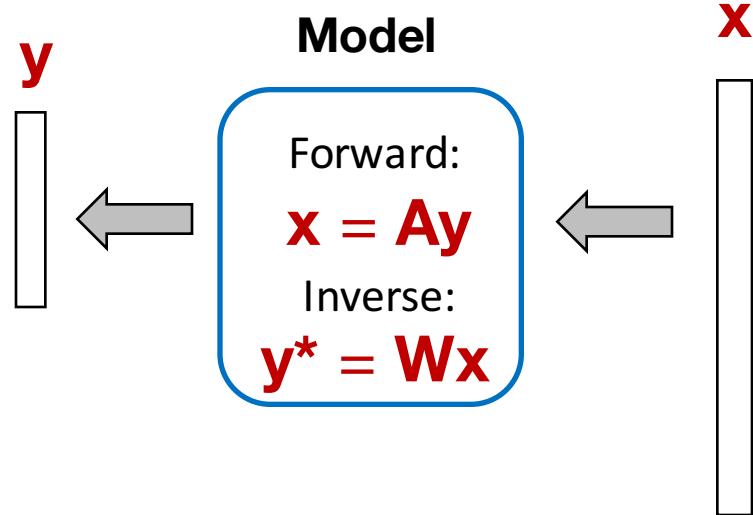
Output: concentrations

Input: spectral measurements

$$\text{Cost function}$$

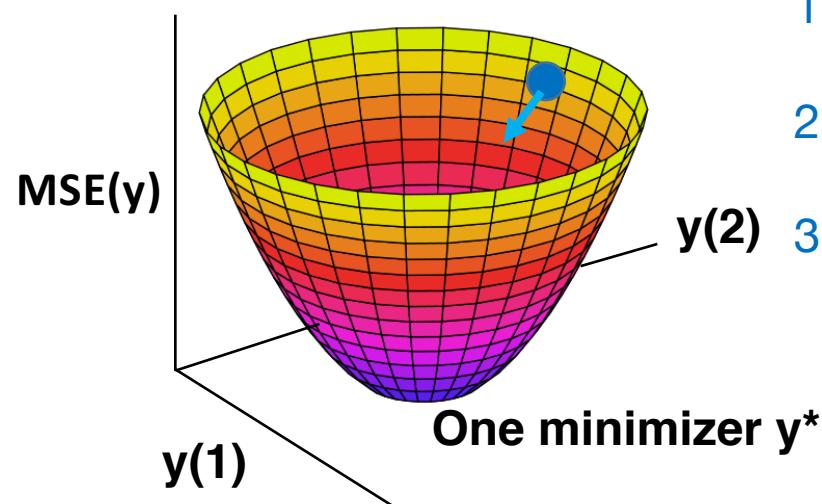
$$0.47$$

$$\text{MSE} = \sum (\mathbf{x} - \mathbf{Ay})^2$$



$$\text{MSE}(\mathbf{y}) = \sum \text{spectral measurements} (\mathbf{x} - \mathbf{Ay})^2$$

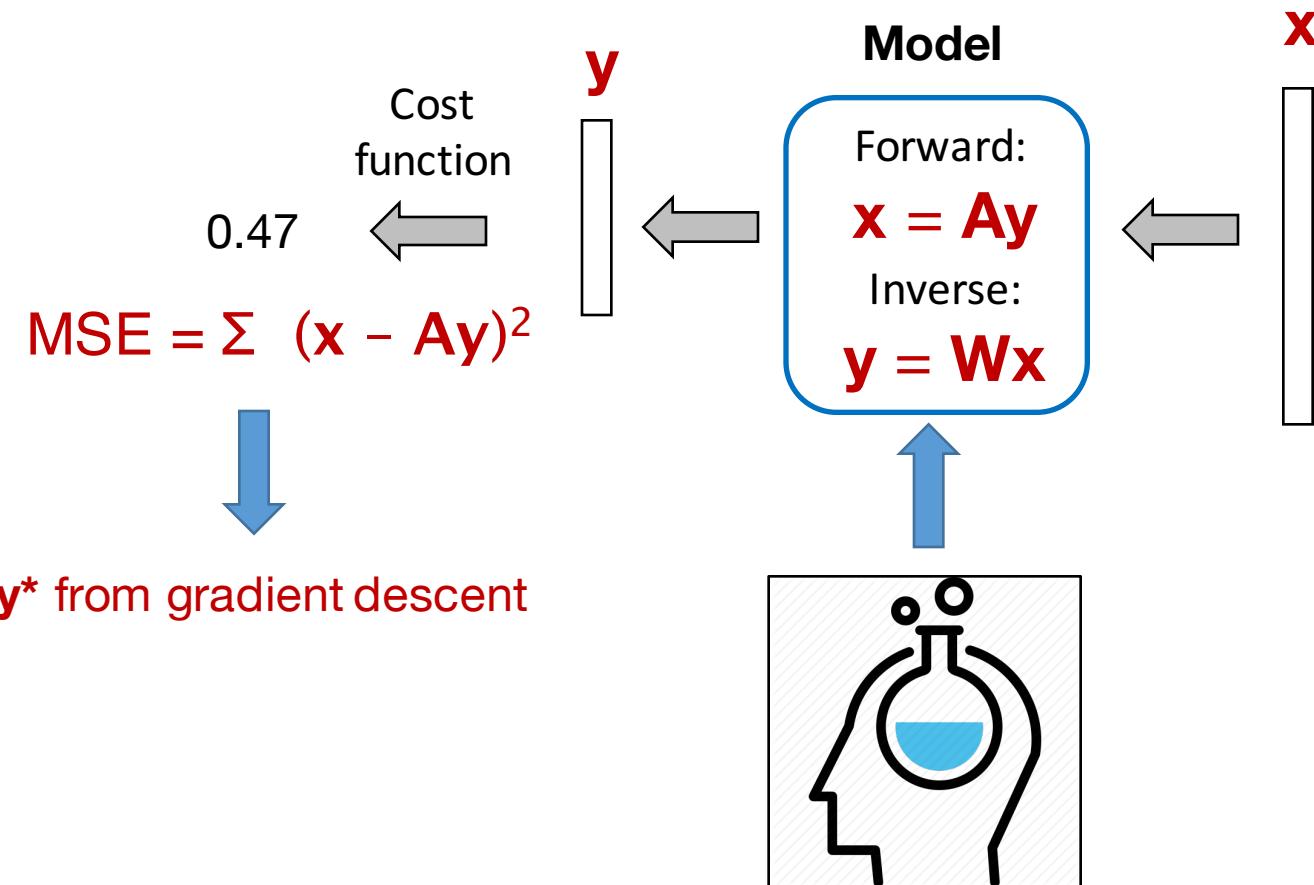
$$\frac{d}{dy} \text{MSE} = \mathbf{A}^T(\mathbf{x} - \mathbf{Ay})$$



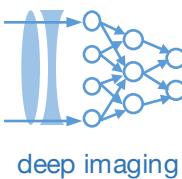
1. Guess output \mathbf{y}
 2. Evaluate slope
 3. If its large, take a step:
- $$\mathbf{y} \leftarrow \mathbf{y} + \varepsilon \mathbf{A}^T(\mathbf{x} - \mathbf{Ay})$$

Optimization pipeline for spectral unmixing

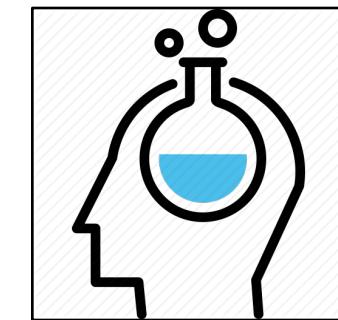
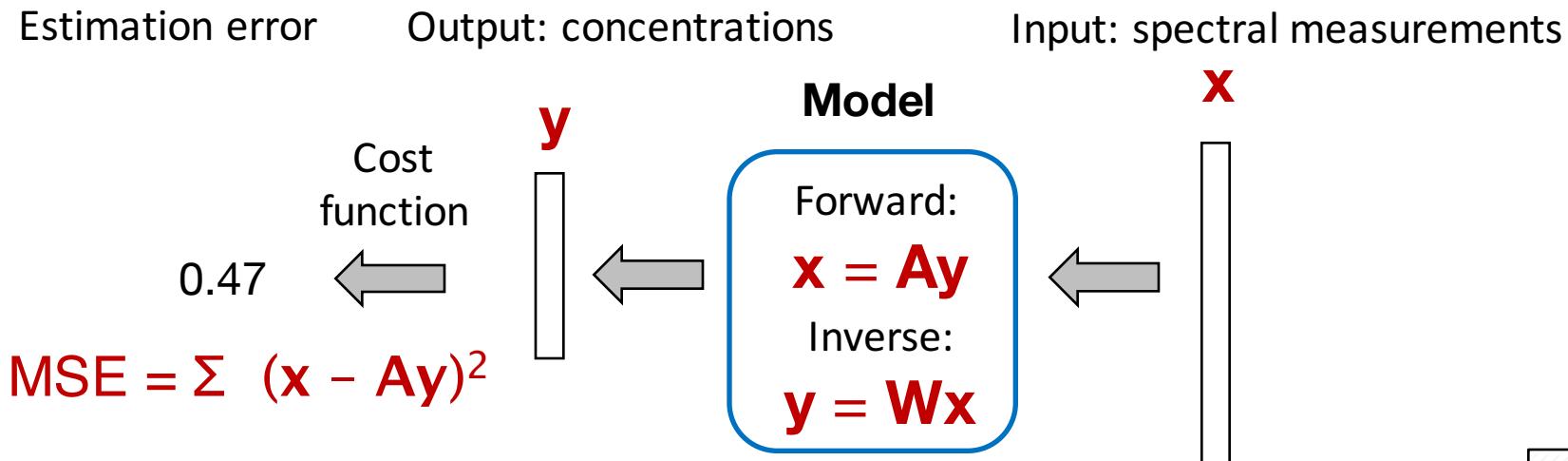
Estimation error Output: concentrations Input: spectral measurements



You or I or someone constructs
A and **W** from first principles
(or something more complex)

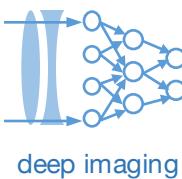


Optimization pipeline for spectral unmixing

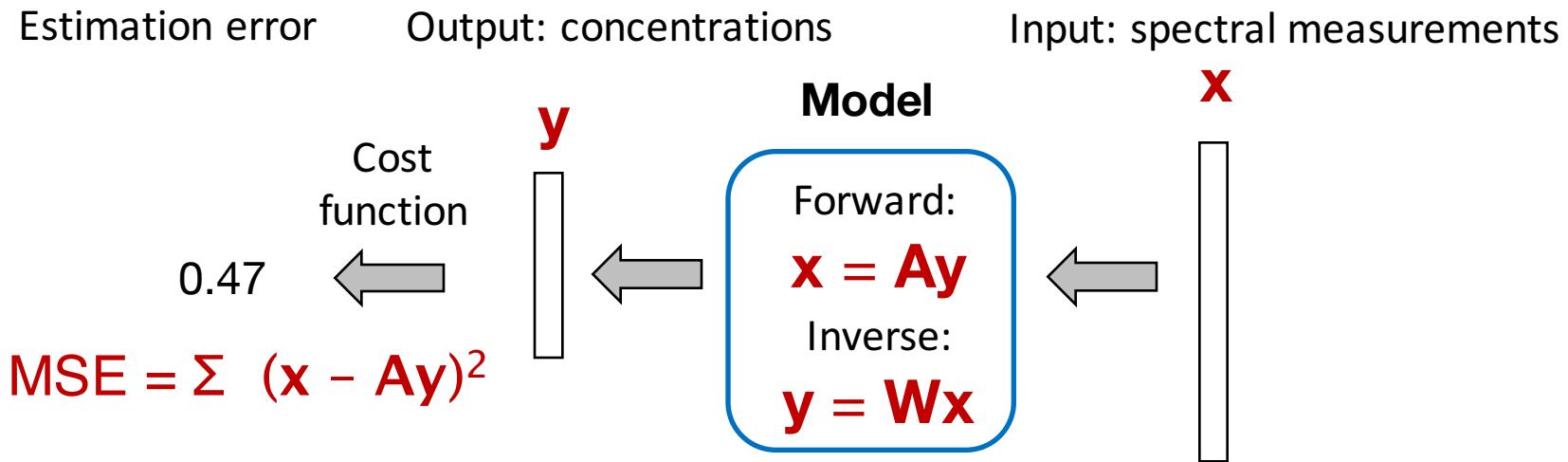


Optimization: You only care about finding the best solution y^*

A and W from first principles

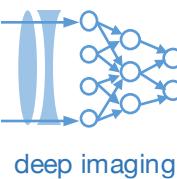


Optimization pipeline for spectral unmixing

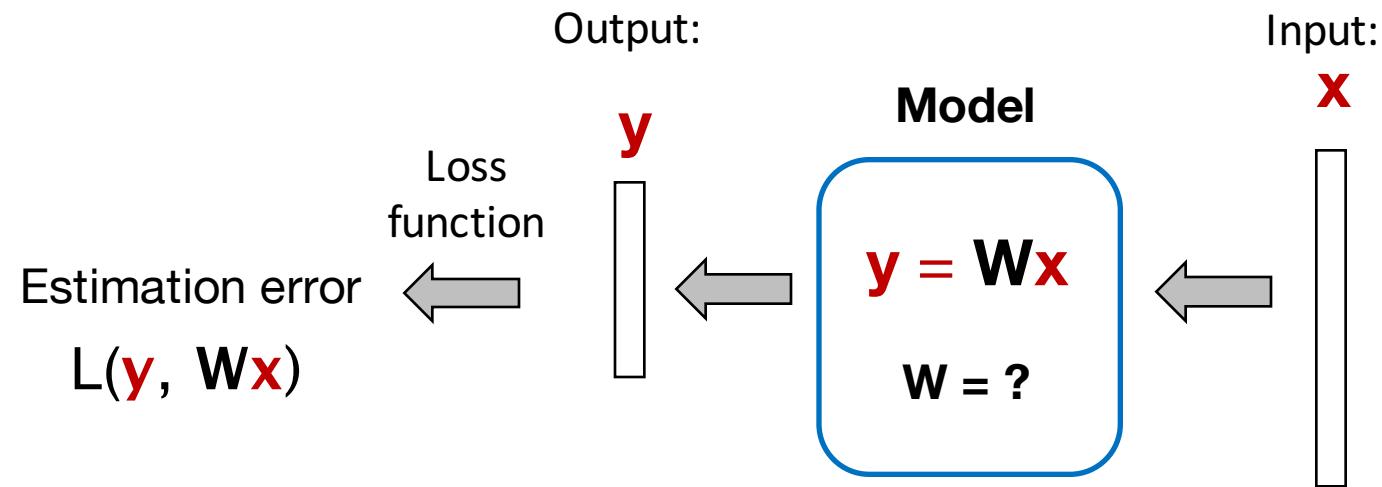


Optimization: You only care about finding the best solution y^*

Machine Learning: You first care about finding the model W , then you'll use that to find the best solution y^*

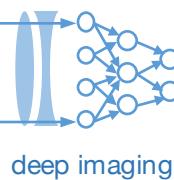


Pipeline for machine learning

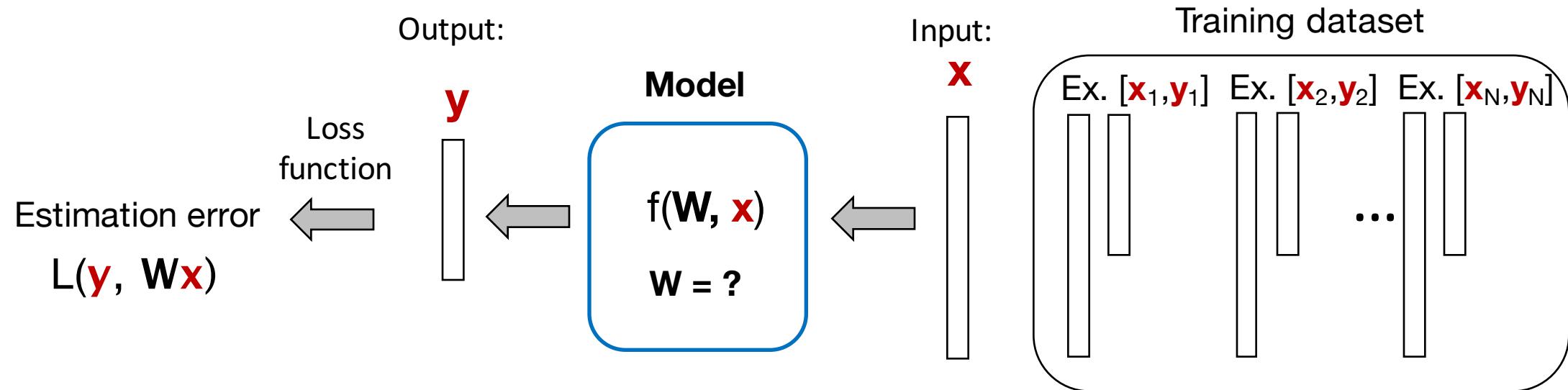


Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix \mathbf{W})

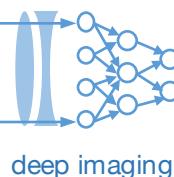


Pipeline for machine learning

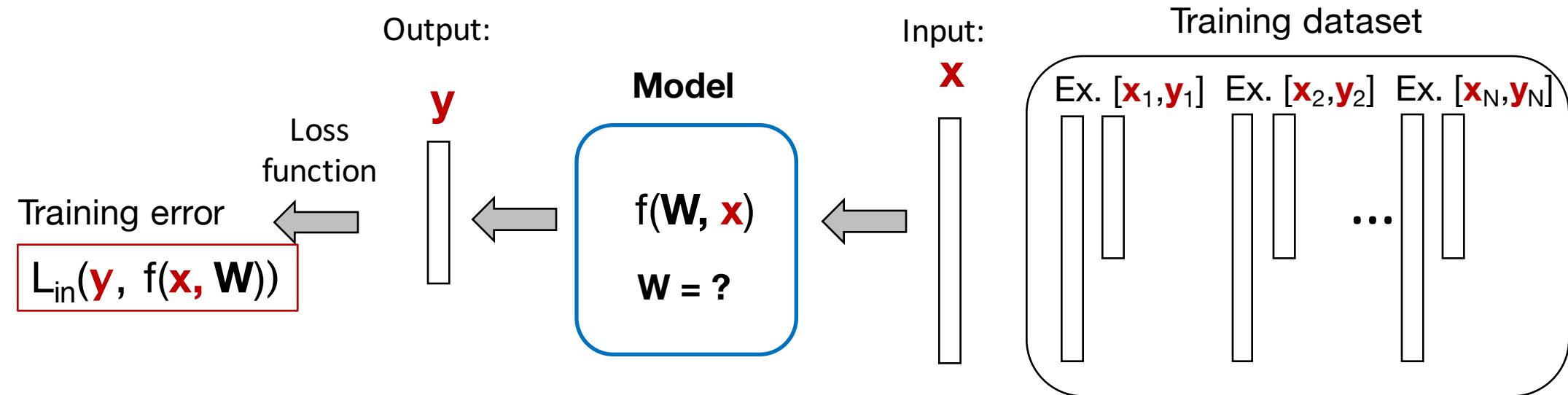


Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix \mathbf{W})
2. Must first determine mapping $f(\mathbf{x}, \mathbf{W})$ using large set of “training” data

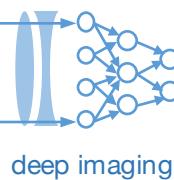


Pipeline for machine learning

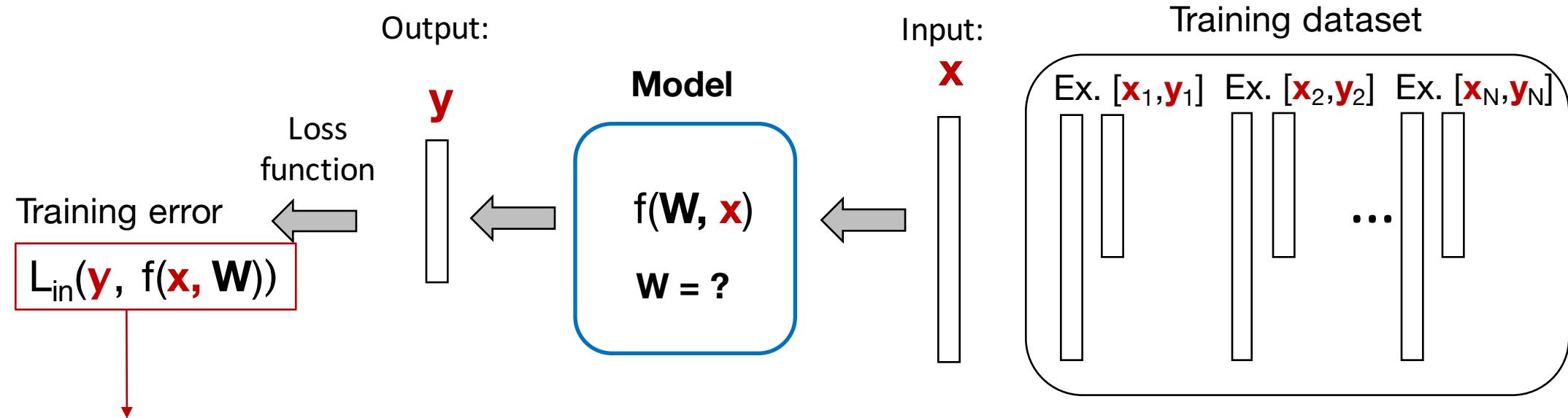


Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix W)
2. Must first determine mapping $f(x, W)$ using large set of “training” data
3. To do so, use a *loss function* L that depends upon the training inputs (x, y) and the model (W)



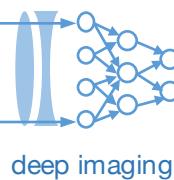
Pipeline for machine learning



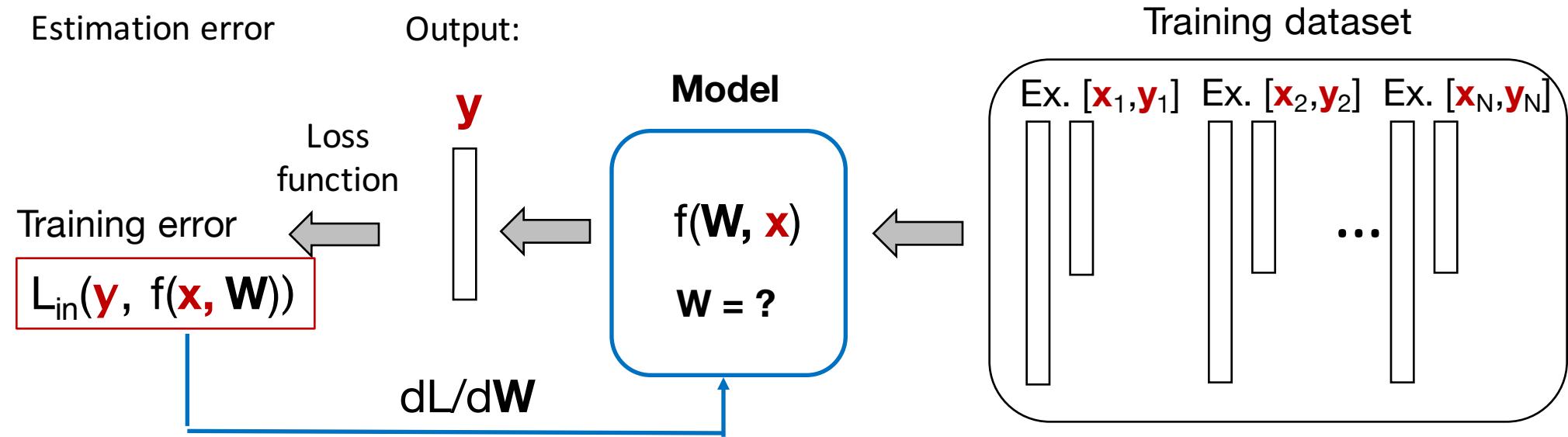
Training Error (“in class error”):

- L_{in} compares modeled output, $f(\mathbf{x}_i, \mathbf{W})$, with the *correct* output that has been *labeled*
- Assume error caused by each labeled example is equally important and sum them up:

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

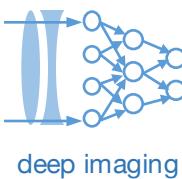


Pipeline for machine learning

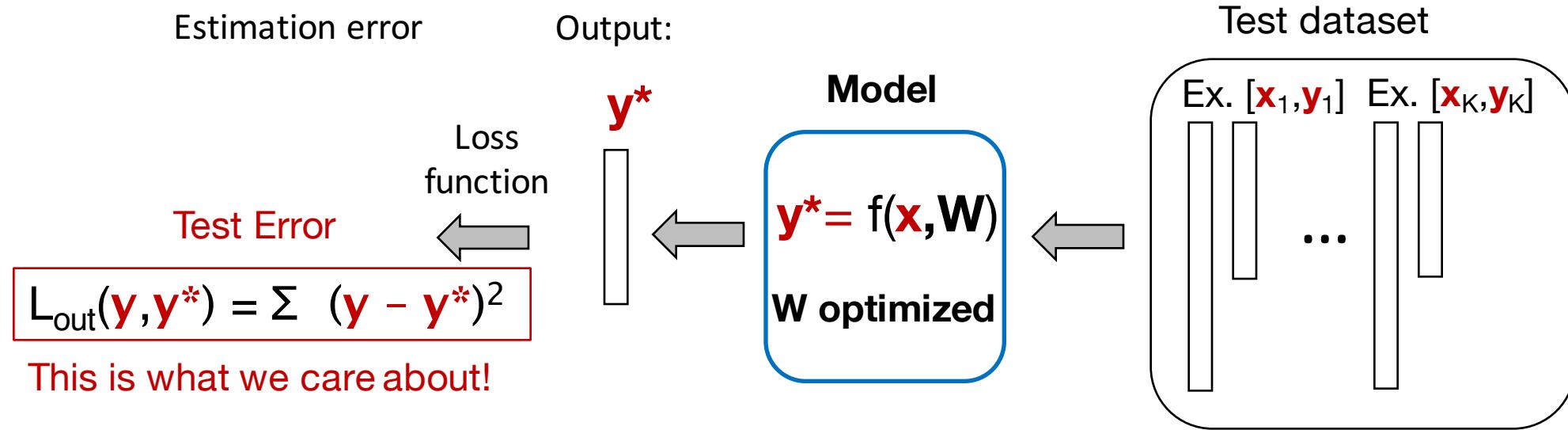


Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix \mathbf{W})
2. Must first determine mapping $f(\mathbf{x}, \mathbf{W})$ using large set of “training” data
3. To do so, use a *loss function* L that depends upon the training inputs (x, y) and the model (\mathbf{W})
4. Find optimal mapping (\mathbf{W}) using the training data, guided by gradient descent on L

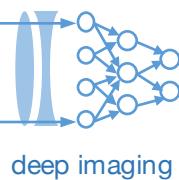


Pipeline for machine learning



Changes for machine learning framework:

1. Now just establish the mapping from inputs to outputs (here, matrix \mathbf{W})
2. Must first determine mapping (\mathbf{W}) using large set of “training” data
3. Use a *loss function* $L(\mathbf{x}, \mathbf{W})$ that depends upon the inputs \mathbf{x} and the model (\mathbf{W})
4. Find optimal mapping (\mathbf{W}), using the training data, guided by gradient descent on L
5. Then, evaluate model accuracy by sending *new* \mathbf{x} through and comparing output y^* to “test” data \mathbf{y}



Example: machine learning for image classification

Correct labels for test images?

Training error

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

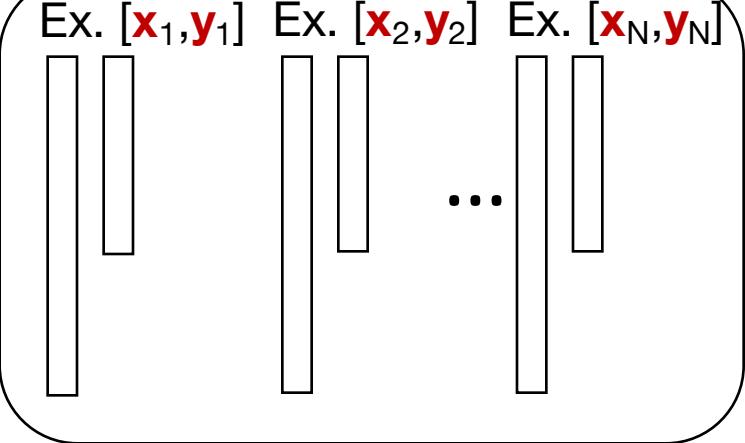
Output: Image class

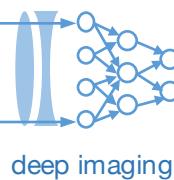
Model

$$\mathbf{y} = f(\mathbf{W}, \mathbf{x})$$

$dL/d\mathbf{W}$

Training data: labeled images





Example: machine learning for image classification

Correct labels for test images?

Training error

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

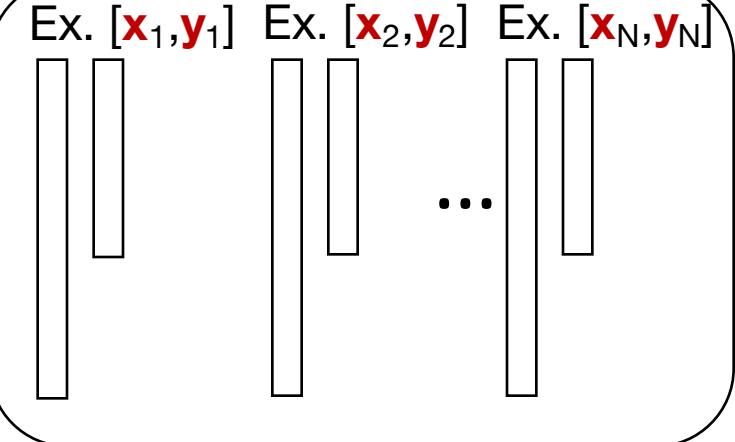
Output: Image class

Model

$$\mathbf{y} = f(\mathbf{W}, \mathbf{x})$$

$dL/d\mathbf{W}$

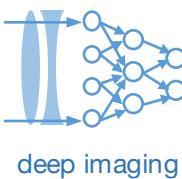
Training data: labeled images



Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

$$\{(x_i, y_i)\}_{i=1}^N$$



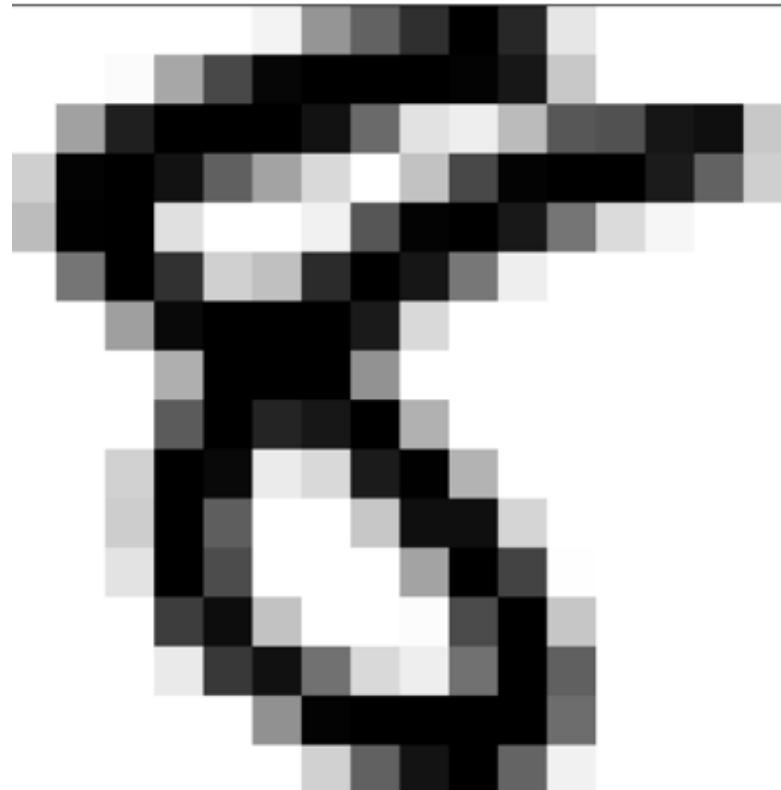
Example: machine learning for image classification

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

https://en.wikipedia.org/wiki/MNIST_database

MNIST image set: <http://yann.lecun.com/exdb/mnist/>

Example: MNIST image dataset



$M = 28 \times 28$ pixel matrix

$m = \text{vec}[M] = 784$ -long vector

Linear model would require $W = 784$ element matrix

Start simple: use $x = (x_0, x_1, x_2)$ to describe **intensity** and **symmetry** of M

Linear model can now use smaller $w = (w_0, w_1, w_2)$

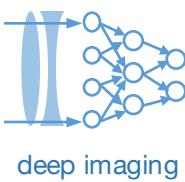
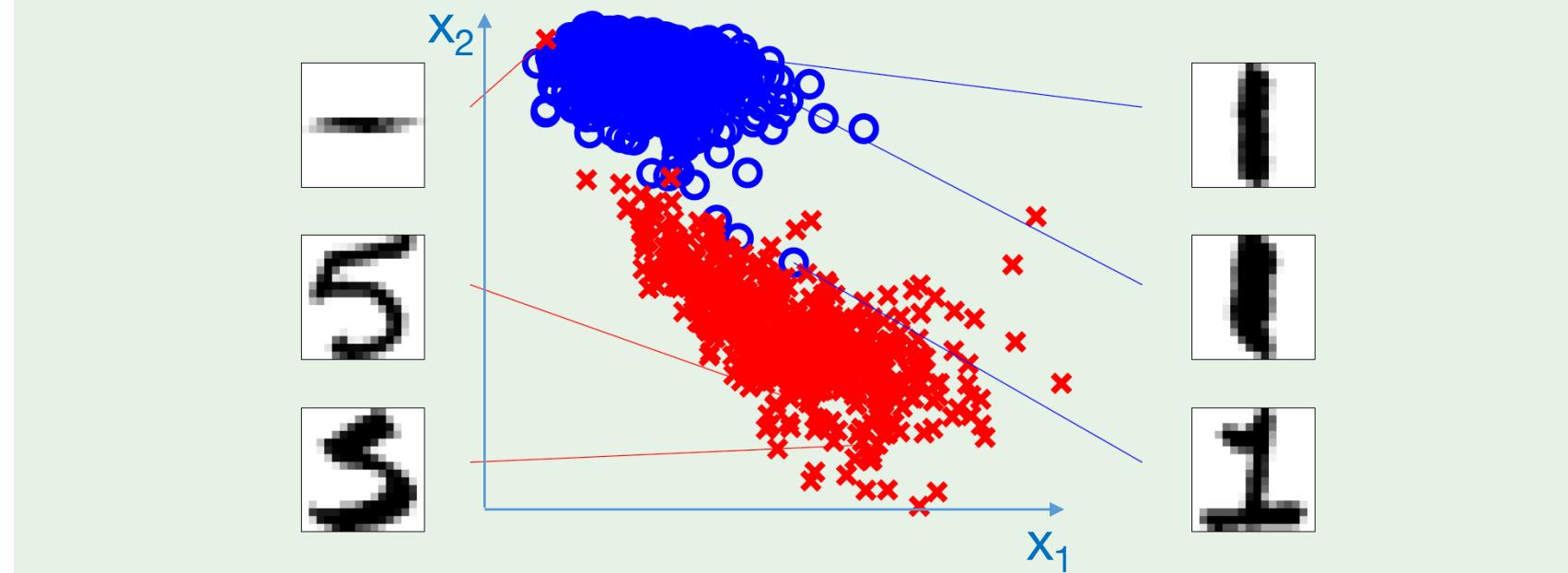


Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

x_1 : intensity

x_2 : symmetry



Caltech Learning from Data: <https://work.caltech.edu/telecourse.html>

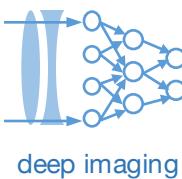
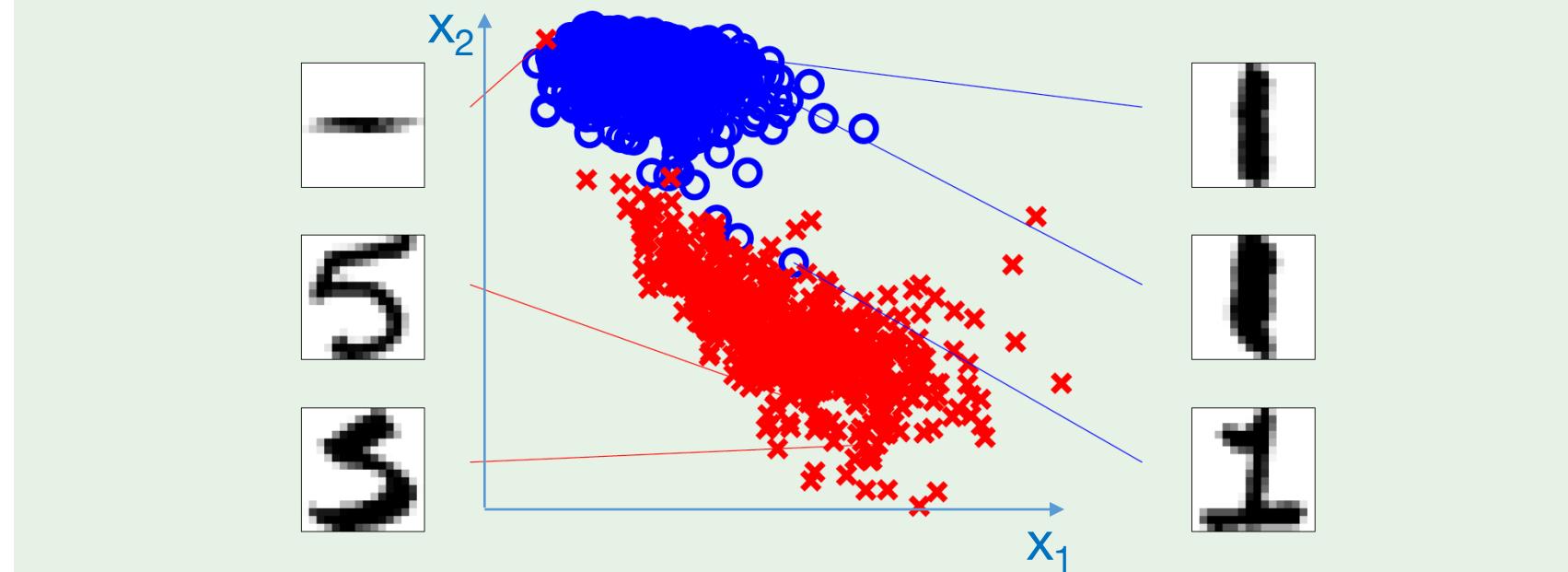


Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

x_1 : intensity

x_2 : symmetry



Dataset: 1000 examples of 1's and 5's mapped to $\mathbf{x}_j = (1, x_1, x_2)$, with associated label $y_j = 1$ or -1

This transforms $wx+b$ into wx

Caltech Learning from Data: <https://work.caltech.edu/telecourse.html>

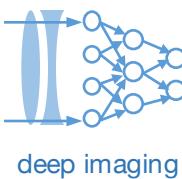
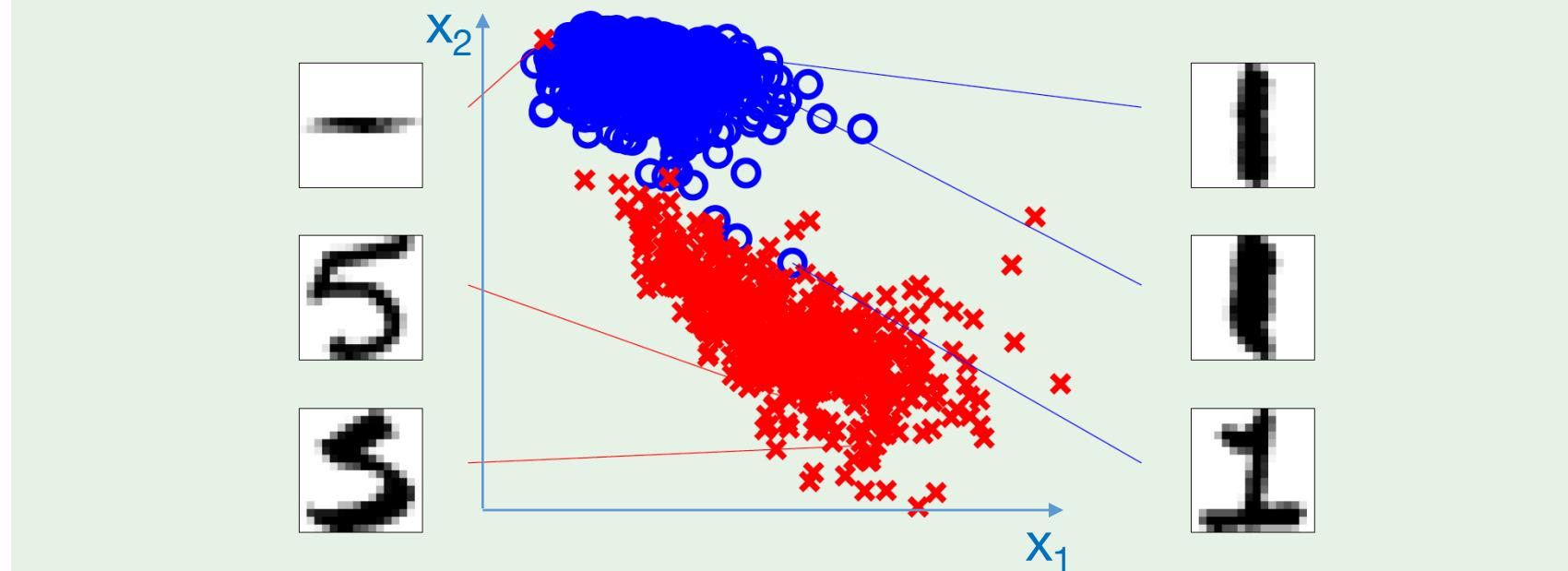


Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

x_1 : intensity

x_2 : symmetry



Dataset: 1000 examples of 1's and 5's mapped to $\mathbf{x}_j = (1, x_1, x_2)$, with associated label $y_j = 1$ or -1

$$[\mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}}] = [\mathbf{x}_j, y_j] \text{ for } n=1 \text{ to } 750$$

$$[\mathbf{X}_{\text{test}}, \mathbf{Y}_{\text{test}}] = [\mathbf{x}_j, y_j] \text{ for } n=751 \text{ to } 1000$$

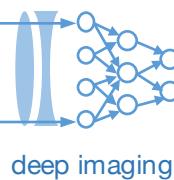
**Labeled
data:**

Training dataset

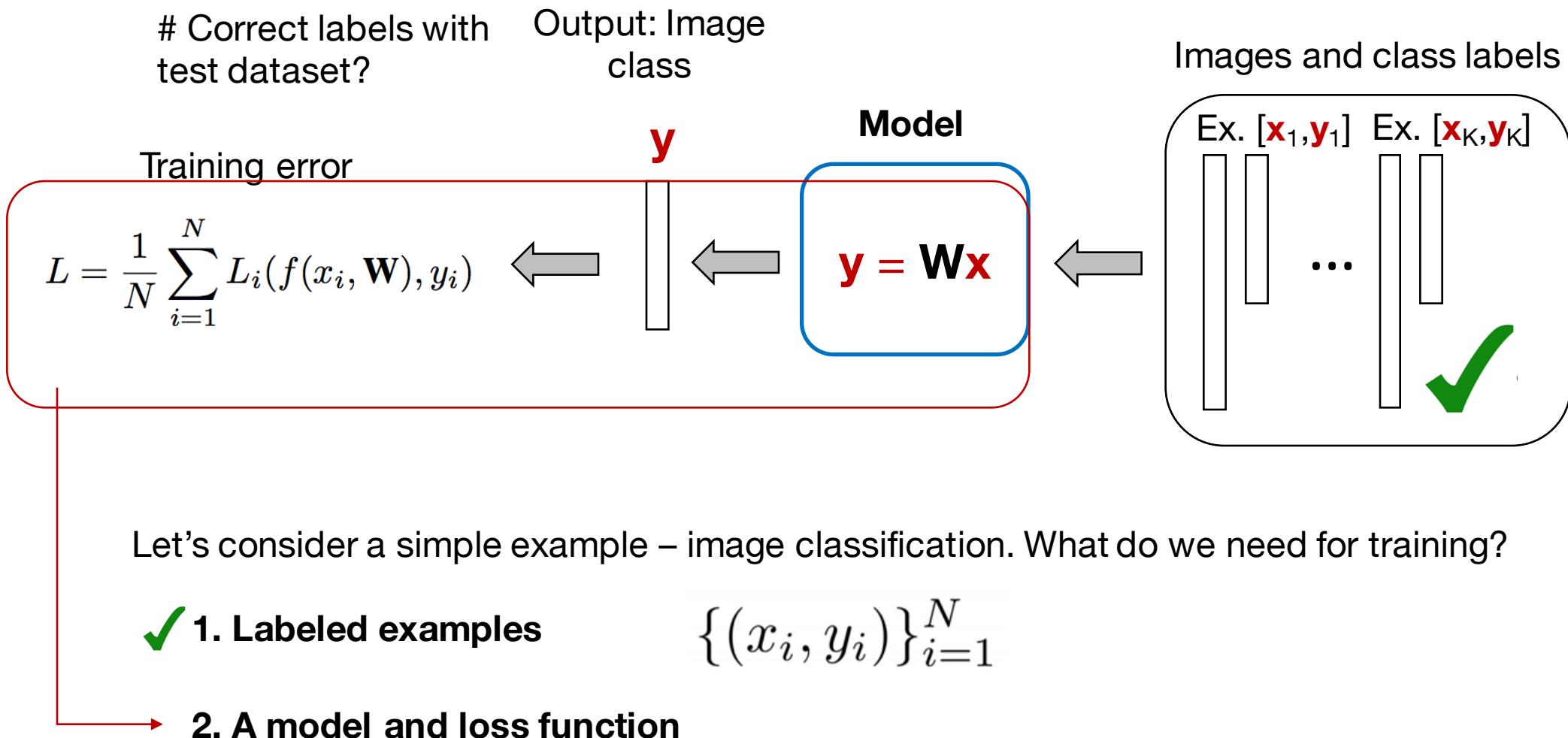
n=1-750

Test dataset

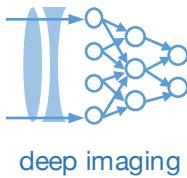
n=751-1000



Example: machine learning for image classification



Let's start with a simpler approach: linear regression



$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

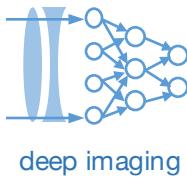
General linear model:

$$L = \frac{1}{N} \sum_{i=1}^N L_i(\mathbf{W}\mathbf{x}_i, y_i)$$

classes \updownarrow

$$\begin{matrix} \mathbf{y} \\ \mathbf{W} \\ \mathbf{x} \end{matrix} = \boxed{\mathbf{W}}$$

Let's start with a simpler approach: linear regression



$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

General linear model:

$$L = \frac{1}{N} \sum_{i=1}^N L_i(\mathbf{W}\mathbf{x}_i, y_i)$$

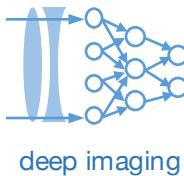
Assume 1 class =
1 linear fit

$$L = \frac{1}{N} \sum_{i=1}^N L_i(w^T x_i - y_i)$$

classes \updownarrow y = $\boxed{\mathbf{w}}$ x

1 var. \updownarrow y = $\boxed{w^T}$ x

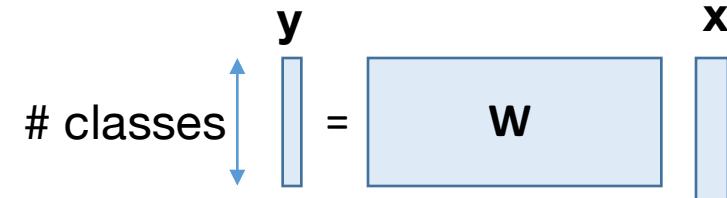
Let's start with a simpler approach: linear regression



$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

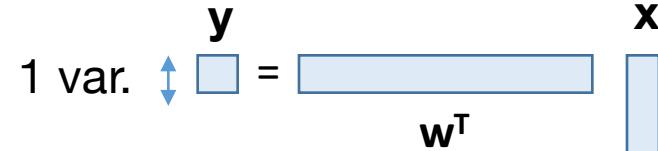
General linear model:

$$L = \frac{1}{N} \sum_{i=1}^N L_i(\mathbf{W}\mathbf{x}_i, y_i)$$



Assume 1 class =
1 linear fit

$$L = \frac{1}{N} \sum_{i=1}^N L_i(w^T x_i - y_i)$$



Use MSE error
model

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

Where labels
determined by
thresholding

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$\text{sgn}(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

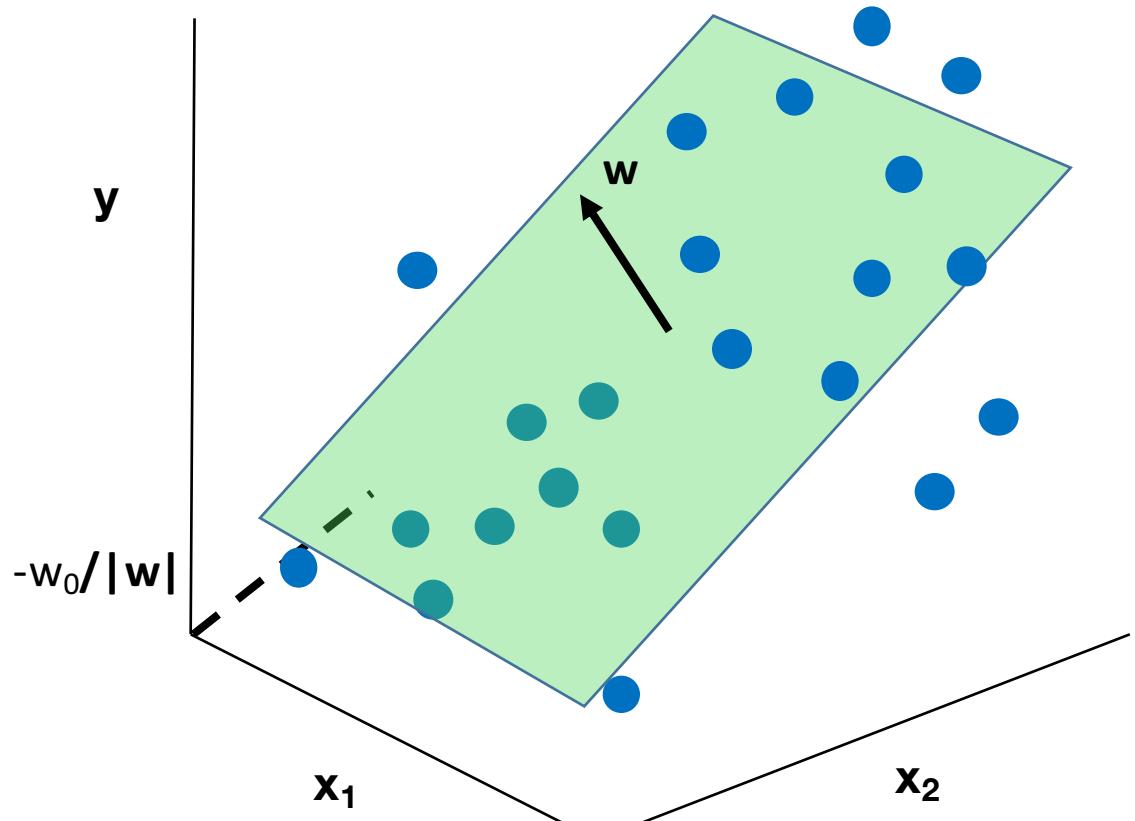
Why does linear regression with $\text{sgn}()$ achieve classification?

Without $\text{sgn}()$: regression for best fit

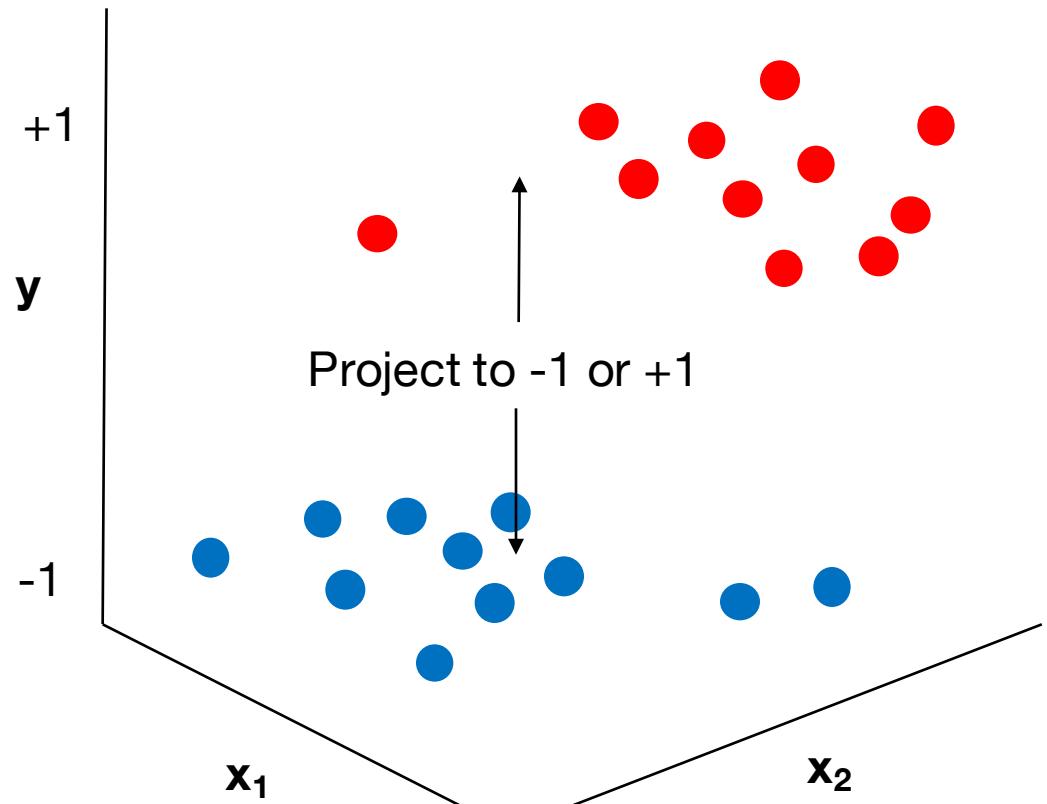
$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- If y_i can be anything, minimizing L makes \mathbf{w} the plane of best fit



Why does linear regression with $\text{sgn}()$ achieve classification?



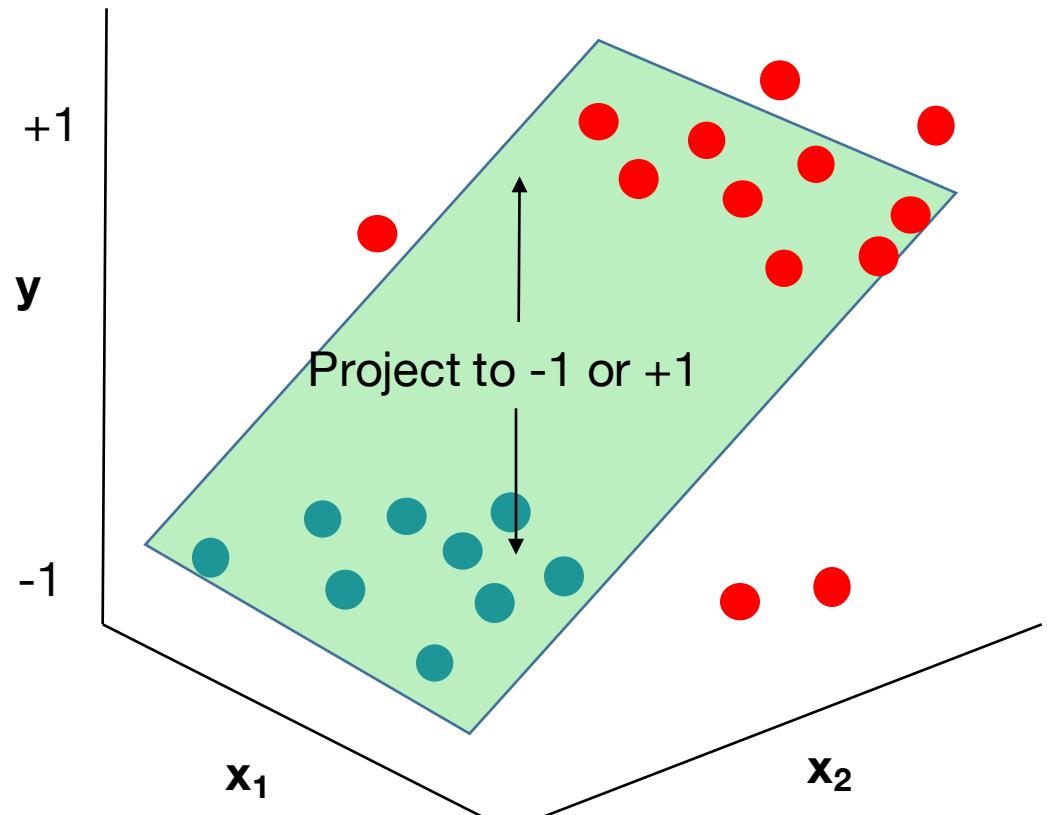
Without $\text{sgn}()$: regression for best fit

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- y_i can only be -1 or $+1$, which defines its class

Why does linear regression with $\text{sgn}()$ achieve classification?



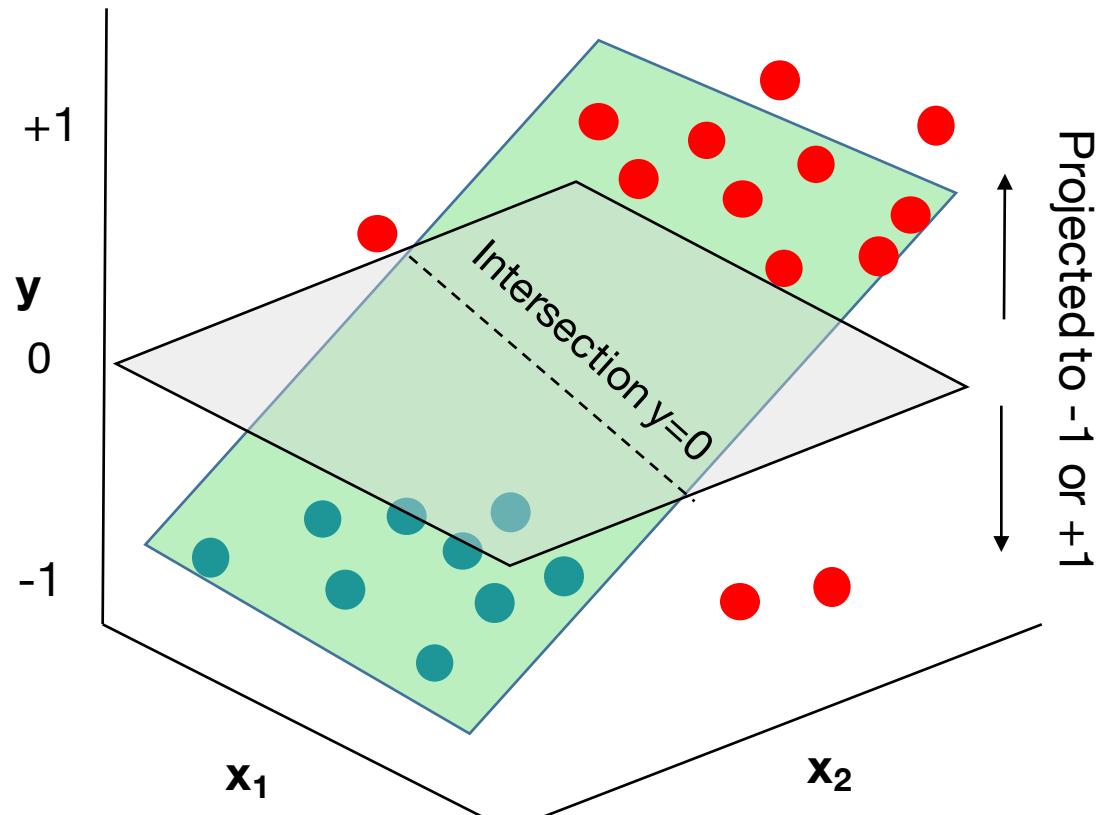
Without $\text{sgn}()$: regression for best fit

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- y_i can only be -1 or $+1$, which defines its class
- Can still find plane of best fit

Why does linear regression with $\text{sgn}()$ achieve classification?



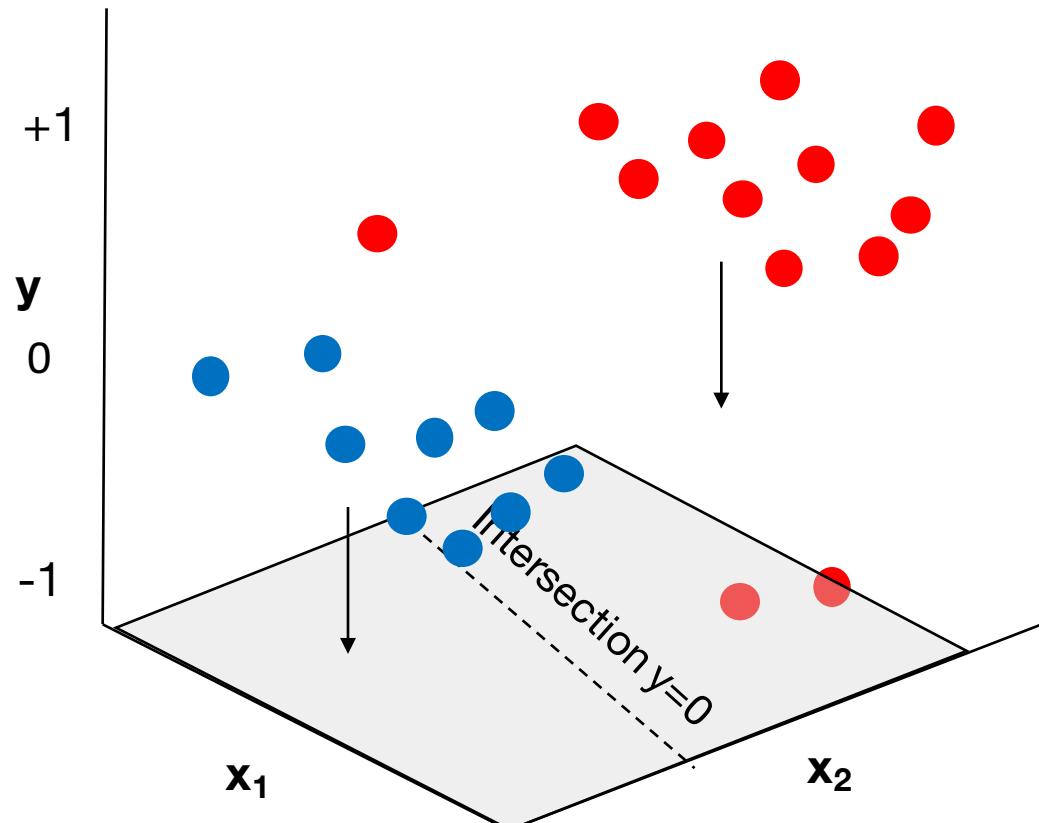
With $\text{sgn}()$ operation:

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- Anything point to one side of $y=0$ intersection is class $+1$, anything on the other side of intersection is class -1

Why does linear regression with $\text{sgn}()$ achieve classification?



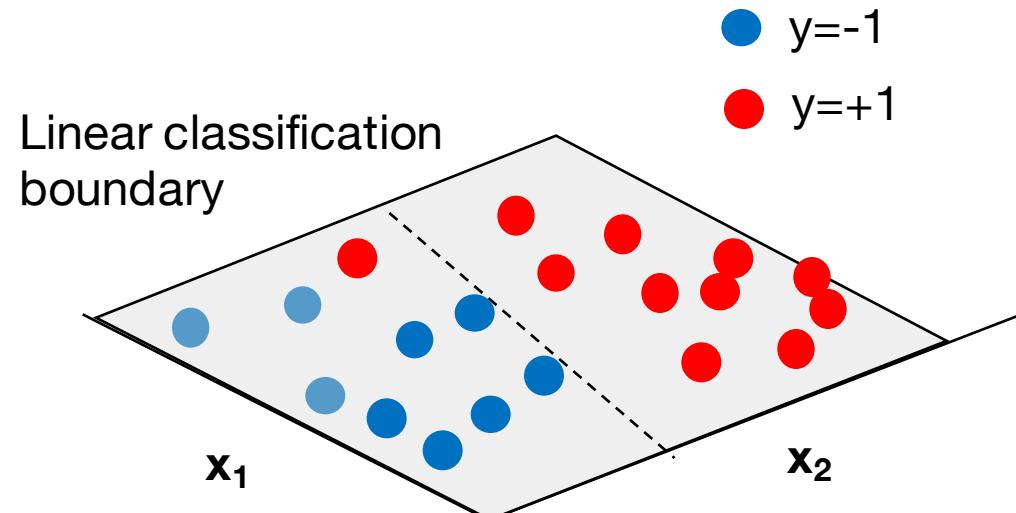
With $\text{sgn}()$ operation:

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- y axis isn't really needed now & can view this decision boundary in 2D

Why does linear regression with $\text{sgn}()$ achieve classification?



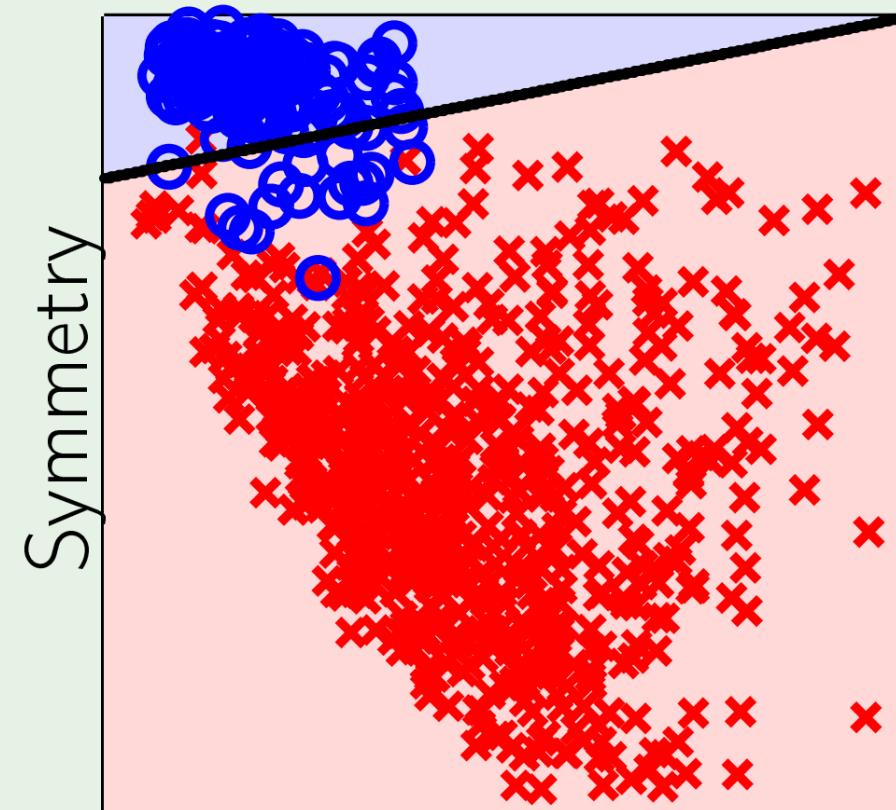
With $\text{sgn}()$ operation:

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

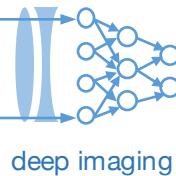
Sign operation takes linear regression and makes it a classification operation!

Linear regression boundary



Average Intensity

Caltech Learning from Data: <https://work.caltech.edu/telecourse.html>



Example: machine learning for image classification

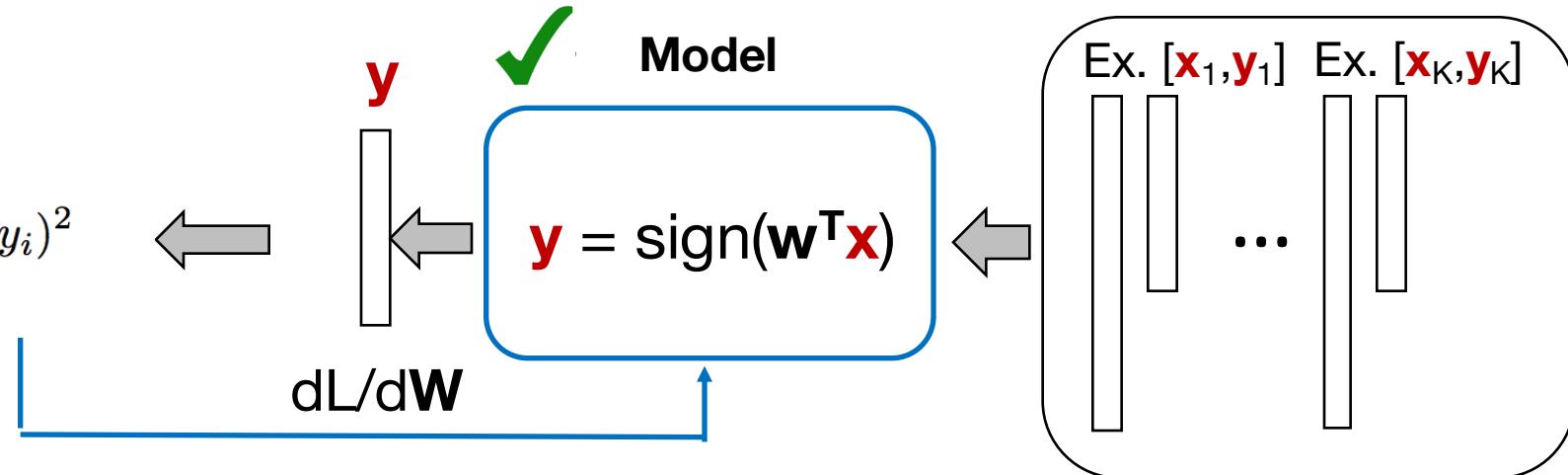


Correct labels with
test dataset?

Output: Image
class

Training error

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

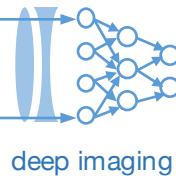


Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

2. A model and loss function

3. A way to minimize the loss function L



Example: machine learning for image classification

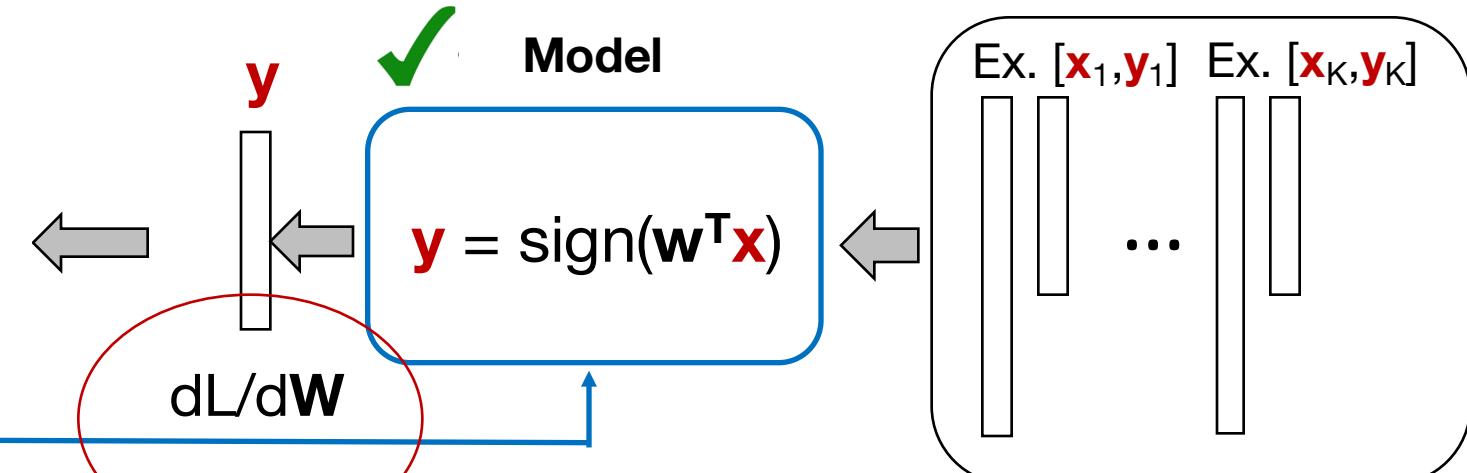


Correct labels with
test dataset?

Output: Image
class

Training error

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

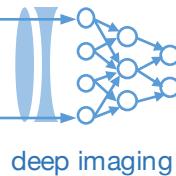


Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

2. A model and loss function

3. A way to minimize the loss function L



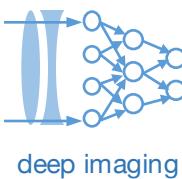
3 methods to solve for w^T in the case of linear regression:

(easier) 1. Pseudo-inverse (this is one of the few cases with a closed-form solution)

2. Numerical gradient descent

3. Gradient descent on the cost function with respect to W

(harder)



1. Turning linear regression for unknown weights W in to a pseudo-inverse:

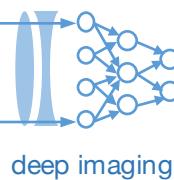
$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

We are multiplying many x_i 's with the same w and are adding them up - let's make a matrix!

$$X = \begin{vmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{vmatrix} \quad \text{Each training image is 1 row of } X$$
$$y = \begin{vmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{vmatrix} \quad \text{Each training label is 1 entry of } y$$

$$L = \frac{1}{N} \|Xw - y\|^2$$

This is the same form as the pseudo-inverse we were working with before, but now we want to solve for w

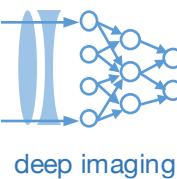


Write this out as a matrix equation:

$$L = \frac{1}{N} \|Xw - y\|^2$$

Note: Training data goes
into “dictionary” matrix

$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + y^T y)$$



Write this out as a matrix equation:

$$L = \frac{1}{N} \|Xw - y\|^2$$

Note: Training data goes into “dictionary” matrix

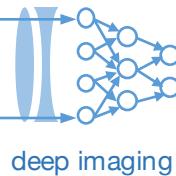
$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + y^T y)$$

$$\nabla L(w) = \frac{2}{N} X^T (Xw - y) = 0$$

Take derivative wrt w and set to 0:

Solution is pseudo-inverse:

$$w_o = (X^T X)^{-1} X^T y$$



Write this out as a matrix equation:

$$L = \frac{1}{N} \|Xw - y\|^2$$

Note: Training data goes into “dictionary” matrix

$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + y^T y)$$

Take derivative wrt w and set to 0:

$$\nabla L(w) = \frac{2}{N} X^T (Xw - y) = 0$$

Solution is pseudo-inverse:

$$w_o = (X^T X)^{-1} X^T y$$

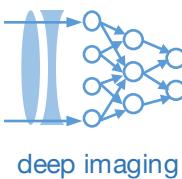
Steps for Pseudo-inverse:

1. Construct matrix X and vector y from data

Each training image is 1 row of X

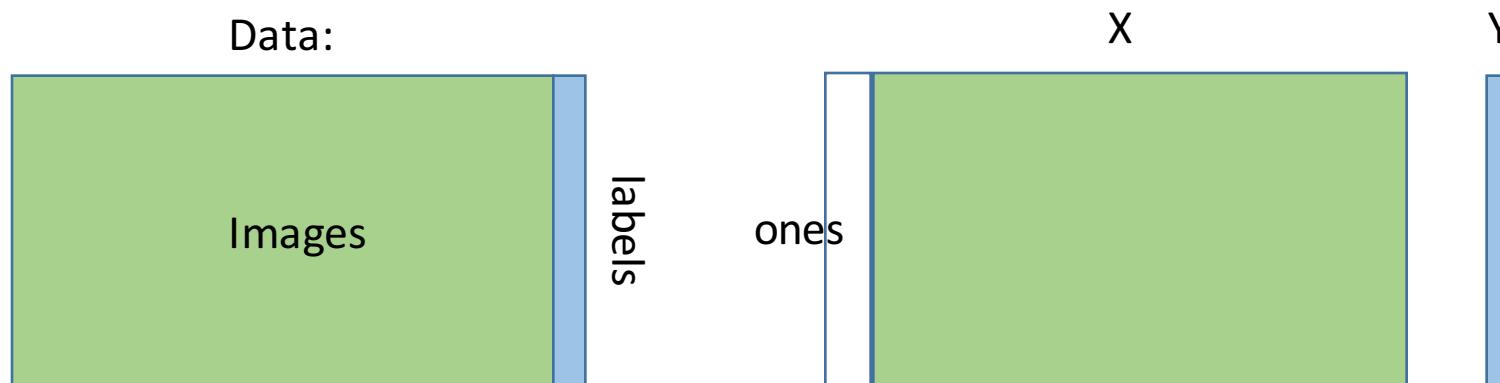
Each training label is 1 entry of y

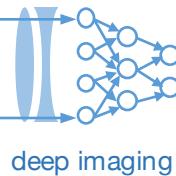
2. Compute solution for w_0 via above equation



Example pseudo-code

```
data = np.loadtxt('train_data.txt', dtype=int)
X = numpy.zeros((data.shape[0],data.shape[1]-1))
X[:,0]=1
Y = numpy.zeros((data.shape[0],1))
for row in m:
    X[row,1:X.shape[1]-1] = data[row,0:data.shape[1]:1]
    Y[row] = data[row,data.shape[1]-1]
X_dagger = np.linalg.pinv(X)
w = np.matmul(X_dagger,Y)
```





3 methods to solve for w^T in the case of linear regression:

(easier) 1. Pseudo-inverse (this is one of the few cases with a closed-form solution)

2. Numerical gradient descent

3. Gradient descent on the cost function with respect to W

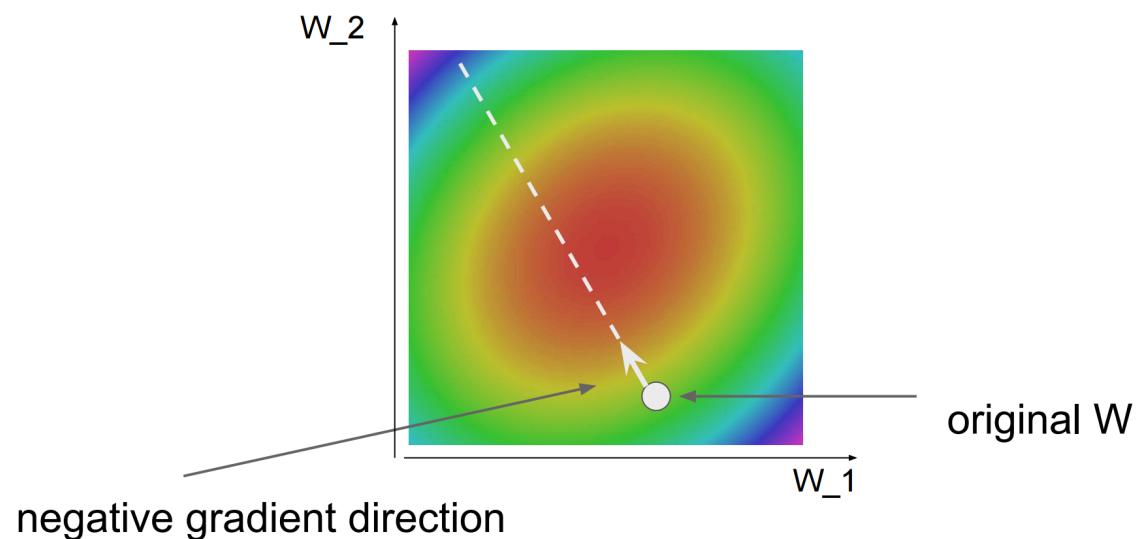
(harder)

Gradient descent: The iterative recipe

Initialize: Start with a guess of \mathbf{W}

Until the gradient does not change very much:
 $dL/d\mathbf{W} = \text{evaluate_gradient}(\mathbf{W}, \mathbf{x}, \mathbf{y}, L)$
 $\mathbf{W} = \mathbf{W} - \text{step_size} * dL/d\mathbf{W}$

`evaluate_gradient` can
be achieved numerically
or algebraically



Gradient descent: Numerical evaluation example



With a matrix, compute this for each entry:

$$\frac{dL(W_i)}{dW_i} = \lim_{h \rightarrow 0} \frac{L(W_i + h) - L(W_i)}{h}$$

Gradient descent: Numerical evaluation example



With a matrix, compute this for each entry:

$$\frac{dL(W_i)}{dW_i} = \lim_{h \rightarrow 0} \frac{L(W_i + h) - L(W_i)}{h}$$

Example:

$$W = [1, 2; 3, 4]$$
$$L(W, x, y) = 12.79$$

Gradient descent: Numerical evaluation example



With a matrix, compute this for each entry:

$$\frac{dL(W_i)}{dW_i} = \lim_{h \rightarrow 0} \frac{L(W_i + h) - L(W_i)}{h}$$

Example:

$$W = [1, 2; 3, 4]$$

$$L(W, x, y) = 12.79$$

$$W_1 + h = [1.001, 2; 3, 4]$$

$$L(W_1 + h, x, y) = 12.8$$

Gradient descent: Numerical evaluation example



With a matrix, compute this for each entry:

$$\frac{dL(W_i)}{dW_i} = \lim_{h \rightarrow 0} \frac{L(W_i + h) - L(W_i)}{h}$$

Example:

$$W = [1, 2; 3, 4]$$

$$L(W, x, y) = 12.79$$

$$W_1 + h = [1.001, 2; 3, 4]$$

$$L(W_1 + h, x, y) = 12.8$$

$$dL(W_1)/dW_1 = 12.8 - 12.79 / .001$$

$$dL(W_1)/dW_1 = 10$$

- Repeat for all entries of W , dL/dW will have $N \times M$ entries for $N \times M$ matrix

3 methods to solve for w^T in the case of linear regression:

(easier) 1. Pseudo-inverse (this is one of the few cases with a closed-form solution)

2. Numerical gradient descent

3. Gradient descent on the cost function with respect to W

(harder)

Next class: We'll go over how gradient descent works, and then examine the cost function to see why it really matters and how we can extend things beyond the simple case of linear regression