

Machine Learning in Imaging

BME 590L
Roarke Horstmeyer

Lecture 5: Ingredients for Machine Learning

Outline

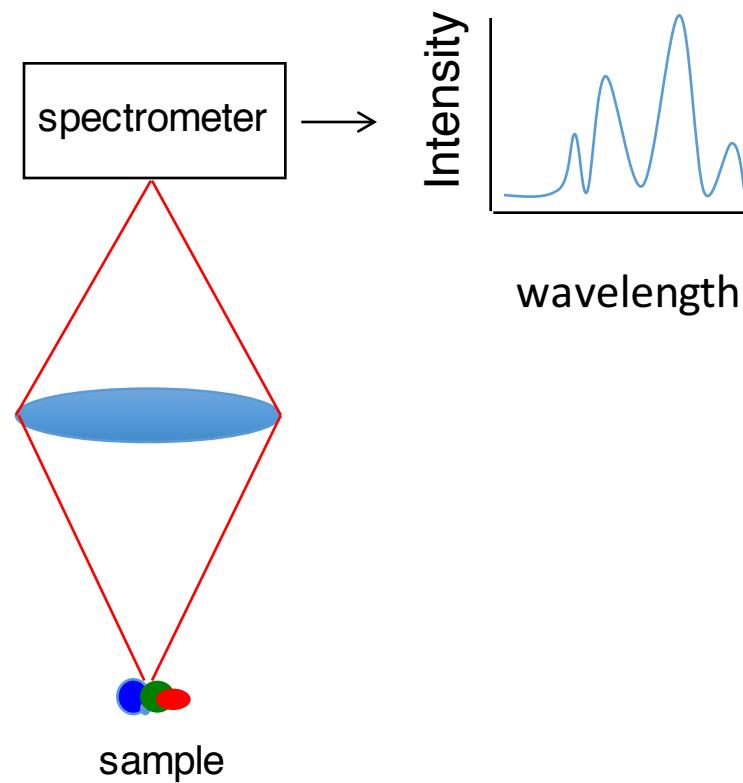
- Review spectral unmixing (last class)
- From optimization to machine learning
- Ingredients for ML
- Example: linear classification of images
 - Train/test data
 - Linear regression model
 - 3 ways to solve

Last time: simple example of spectral unmixing

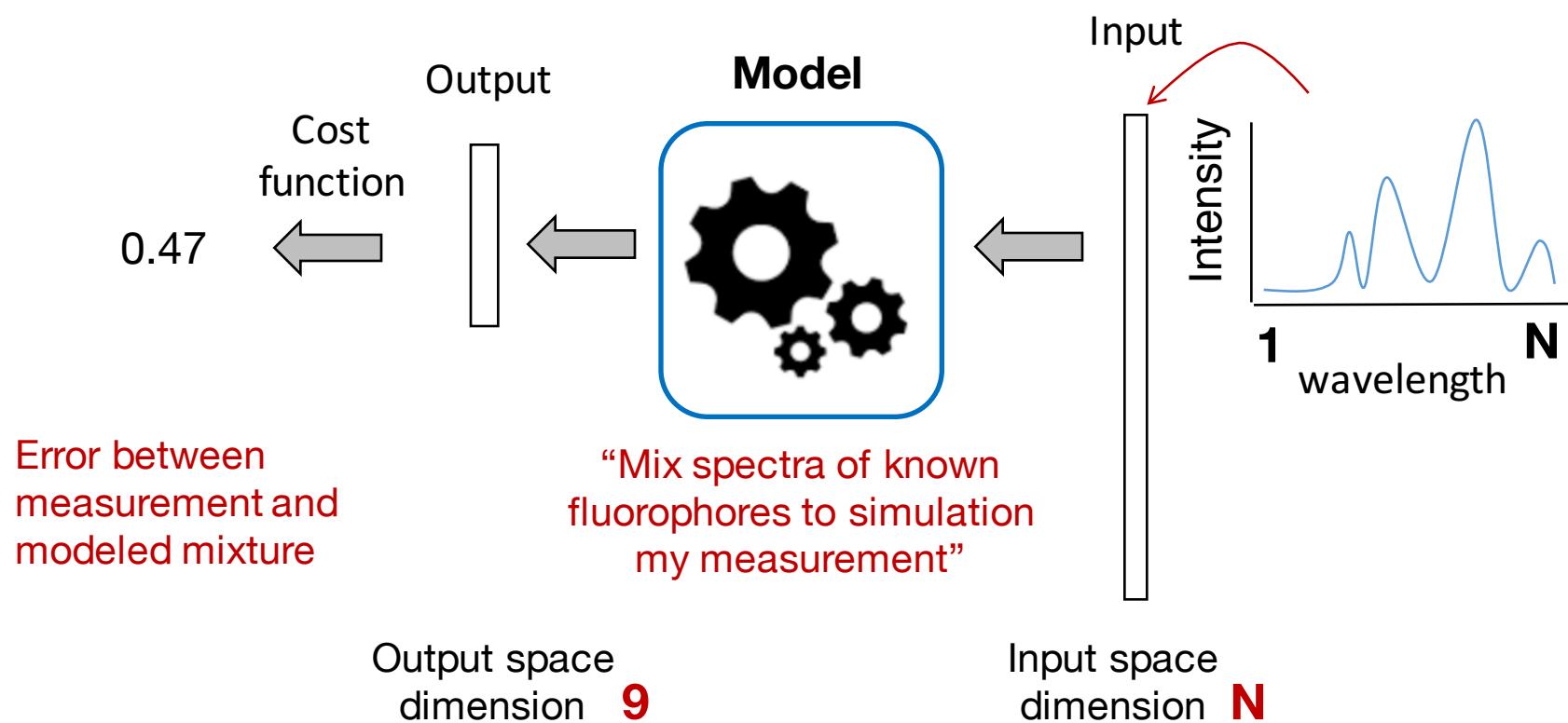
(For whatever reason, whenever I get confused about optimization, I think about this example)

The setup:

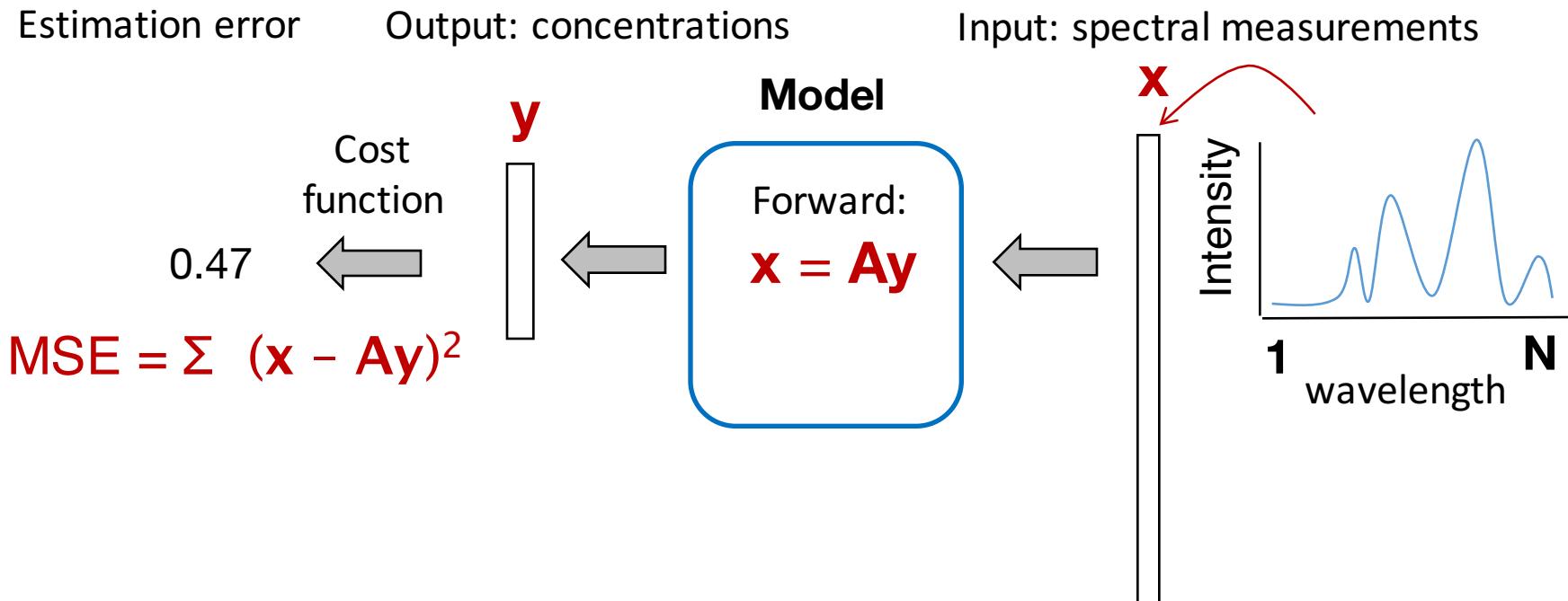
- measure the color (spectral) response of a sample (e.g., how much red, green and blue there is, or several hundred measurements of its different colors).
- You know that the sample can only contain 9 different fluorophores.
- What % of each fluorophores is in your sample?



Optimization pipeline for spectral unmixing

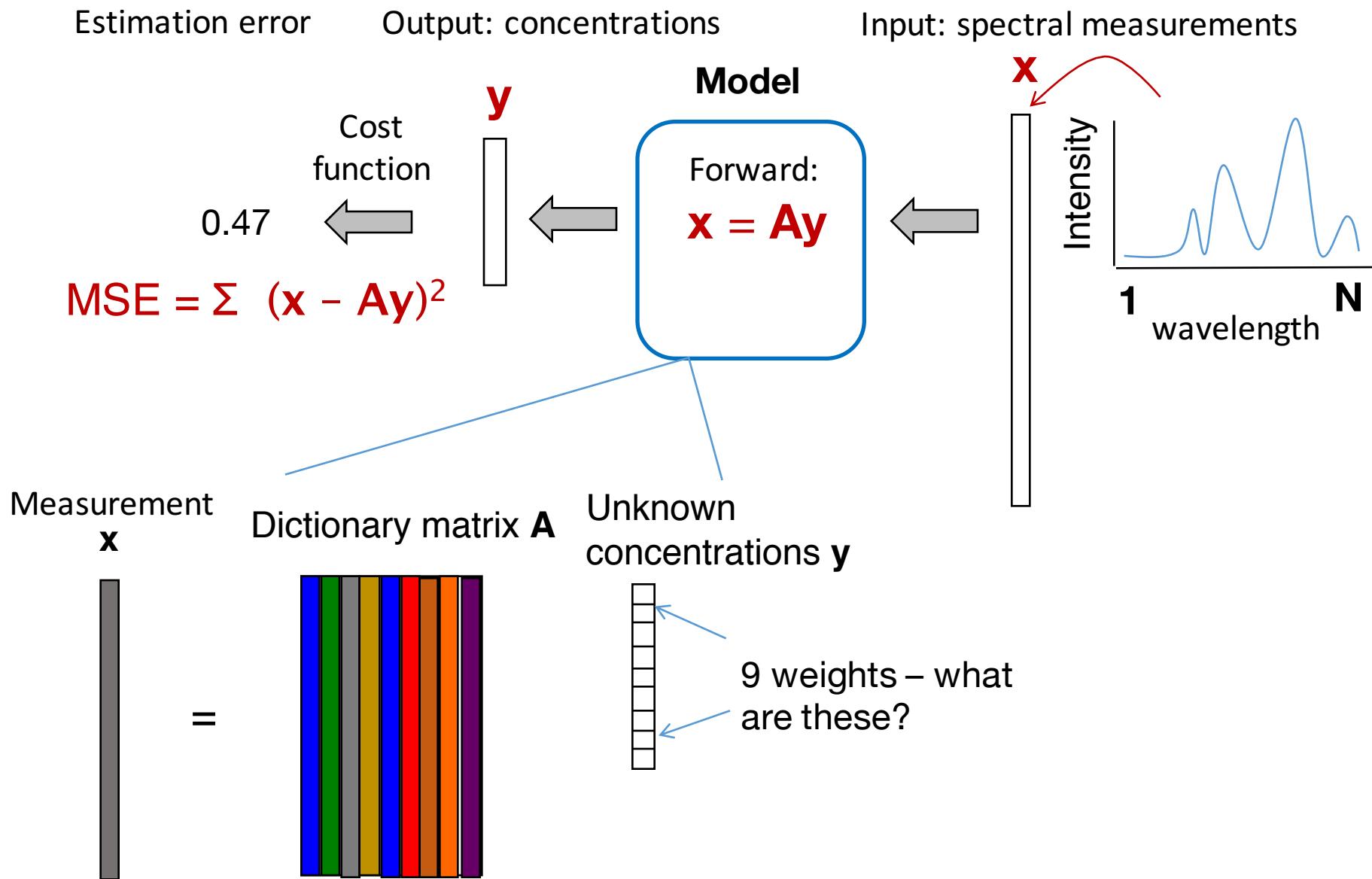


Optimization pipeline for spectral unmixing

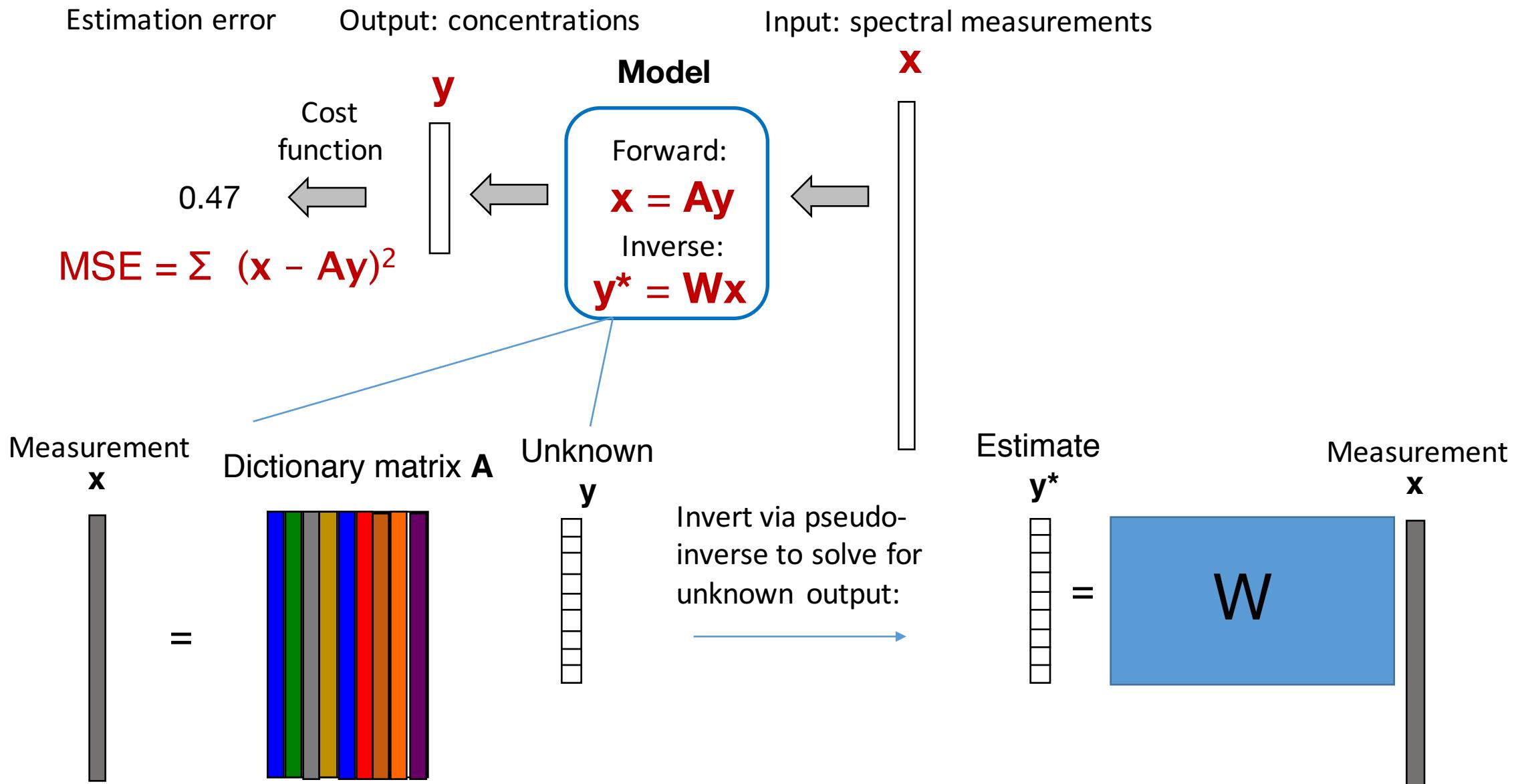


*Note: notation changed from last time to be consistent with we'll use in the future

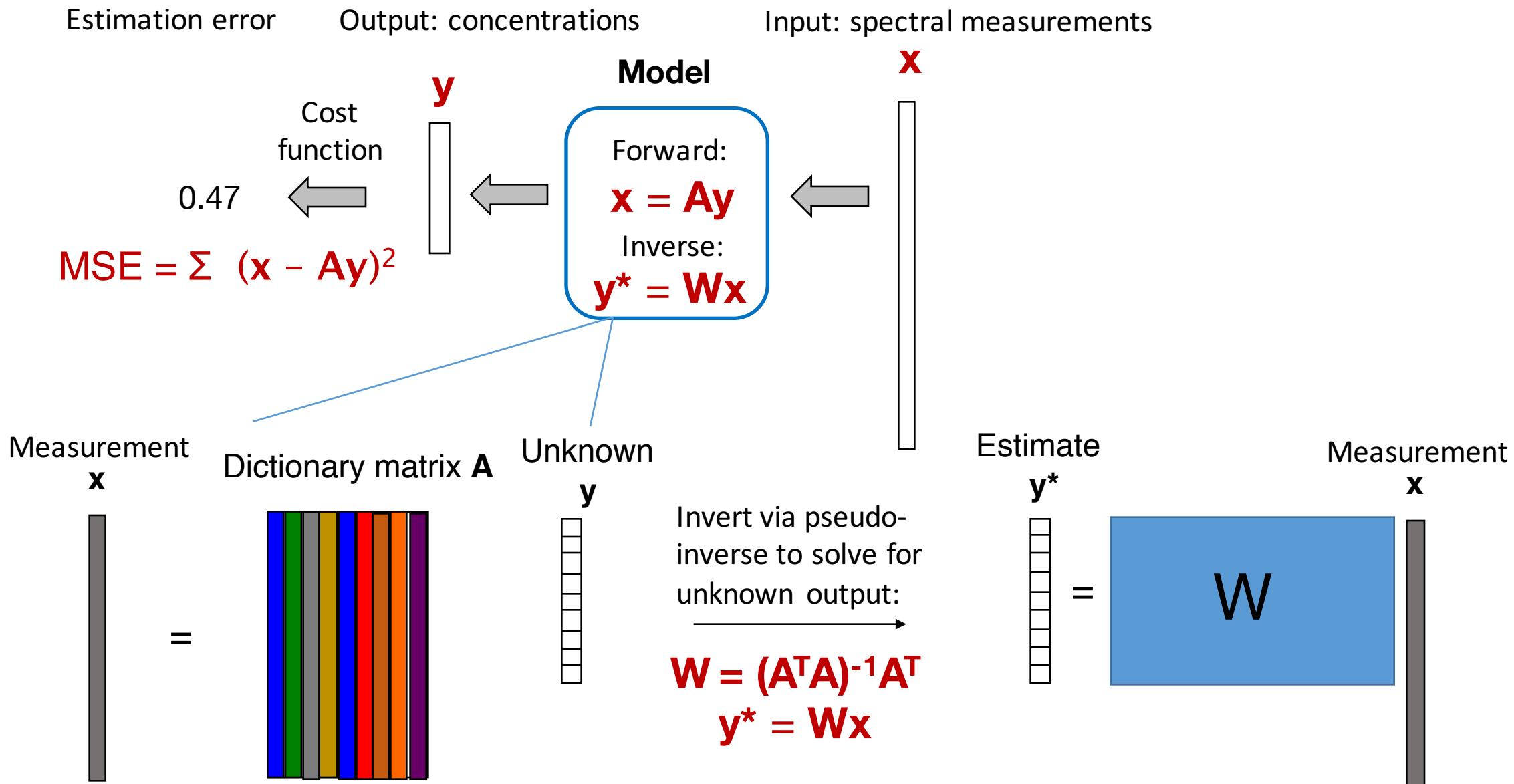
Optimization pipeline for spectral unmixing



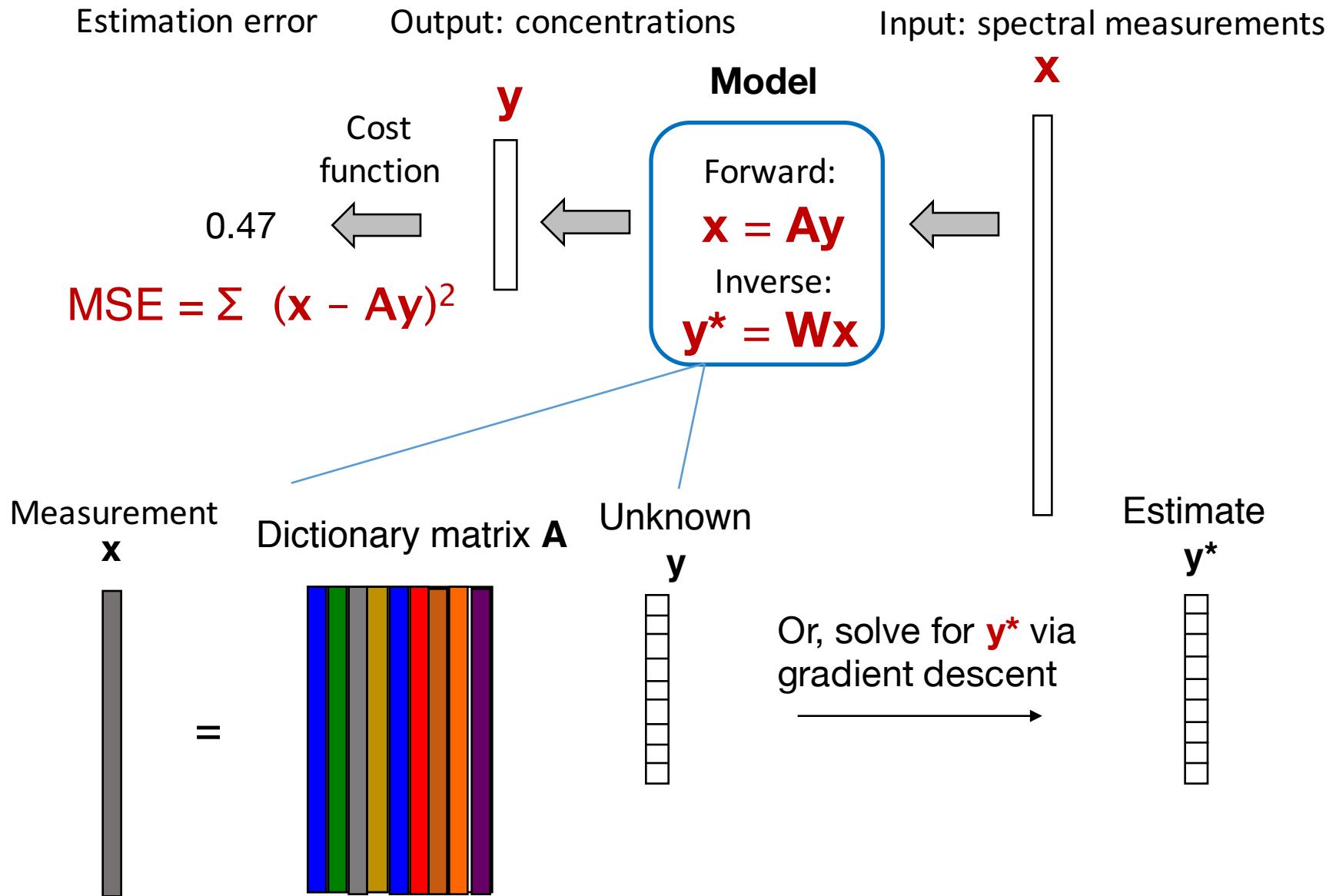
Optimization pipeline for spectral unmixing



Optimization pipeline for spectral unmixing



Optimization pipeline for spectral unmixing

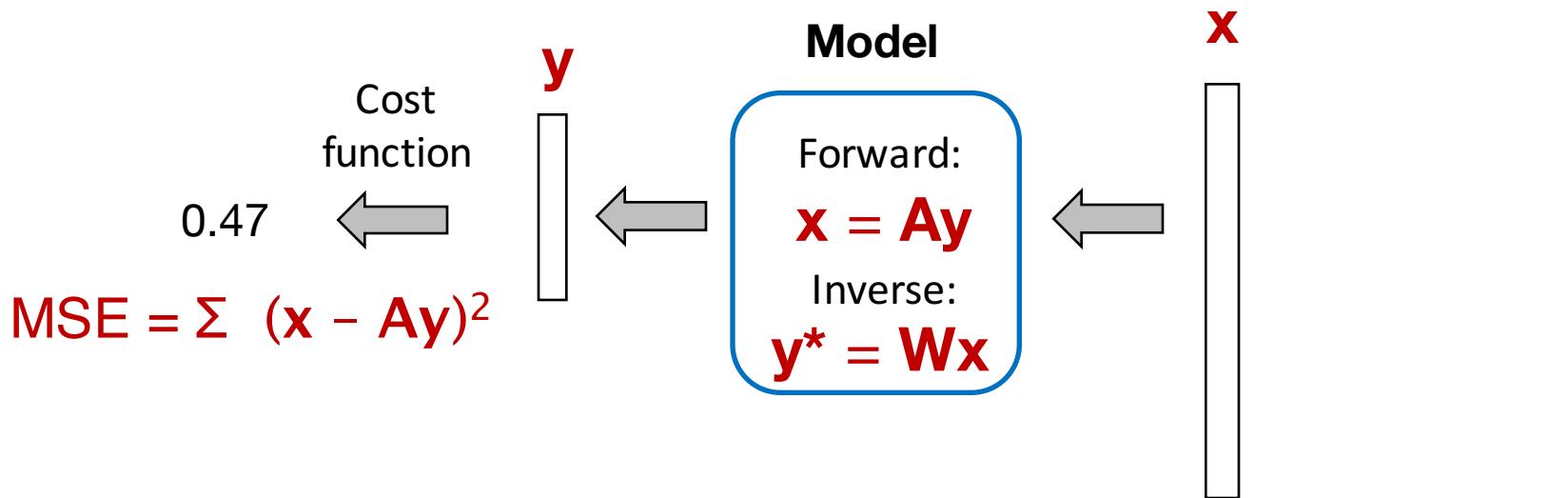


Optimization pipeline for spectral unmixing

Estimation error

Output: concentrations

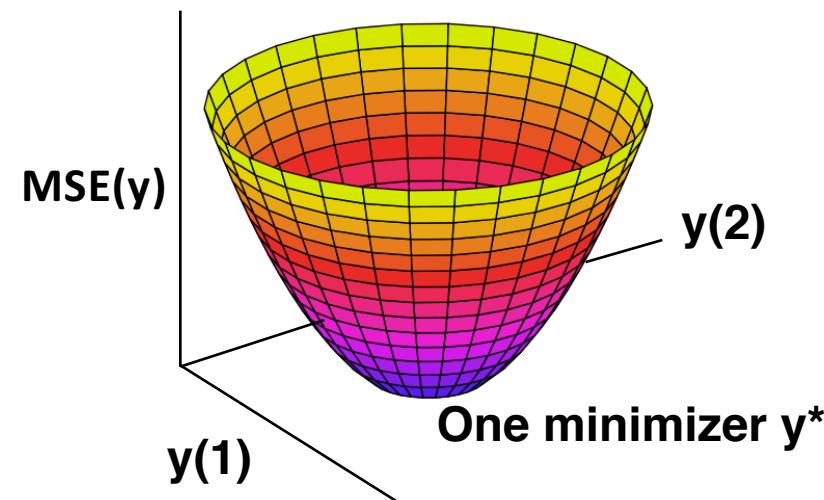
Input: spectral measurements



$$MSE(y) = \sum (x - Ay)^2$$

spectral
measurements

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

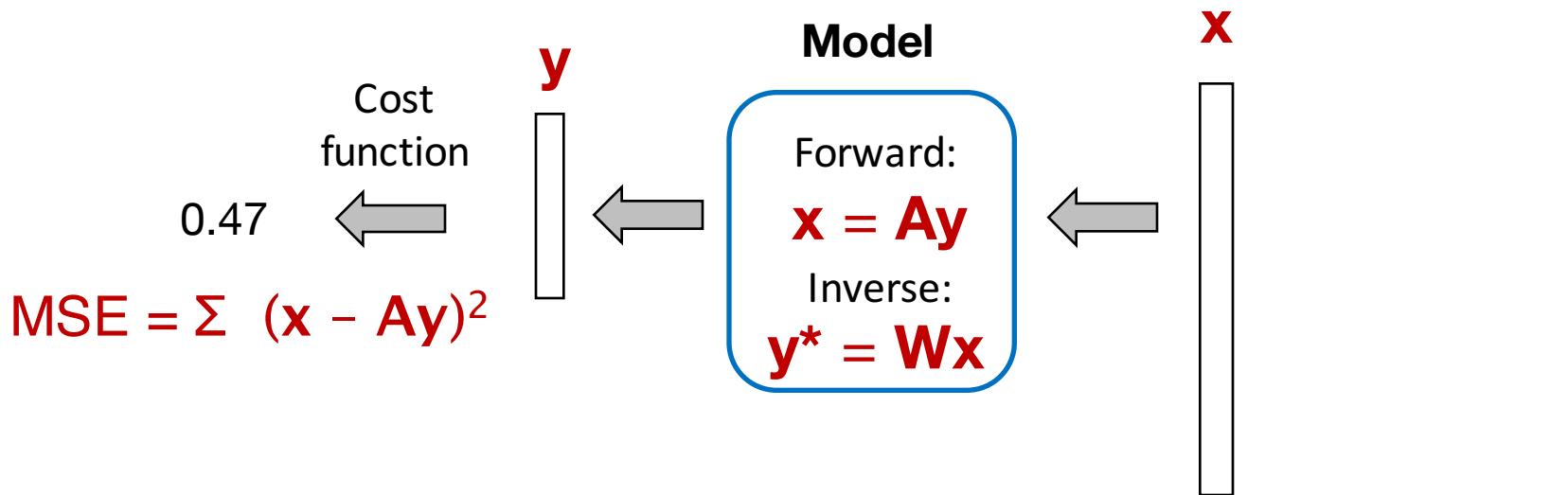


Optimization pipeline for spectral unmixing

Estimation error

Output: concentrations

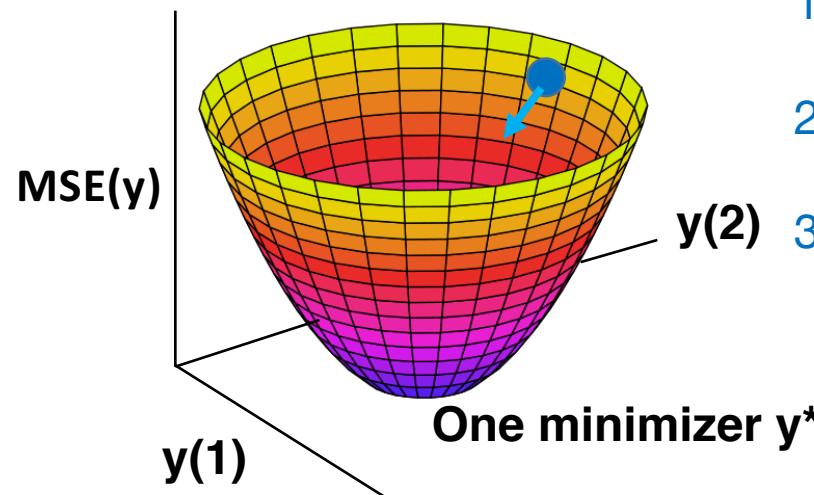
Input: spectral measurements



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

spectral
measurements

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



1. Guess output \mathbf{y}
2. Evaluate slope
3. If its large, take a step:

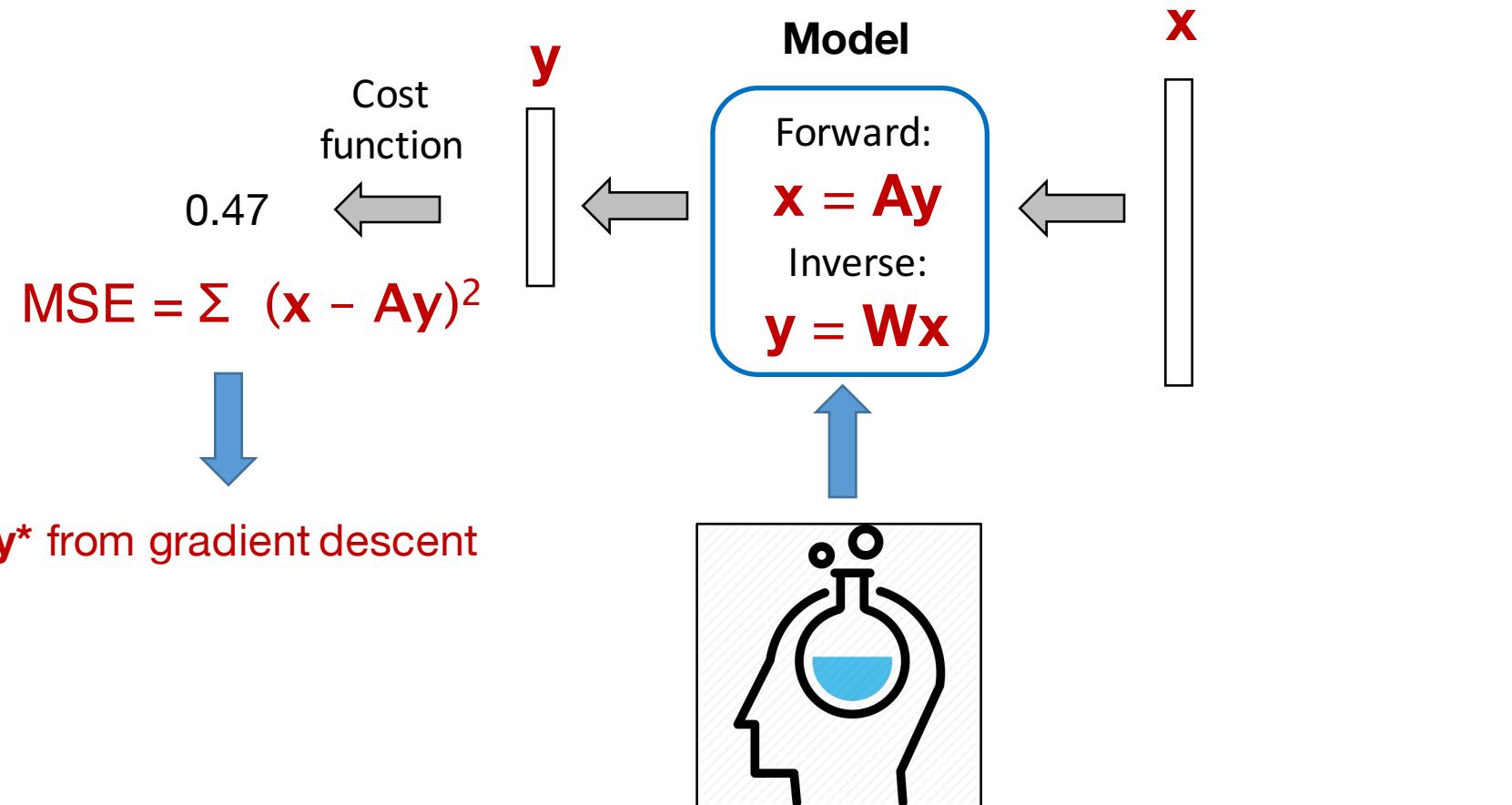
$$\mathbf{y} \leftarrow \mathbf{y} + \varepsilon \mathbf{A}^T(\mathbf{x} - \mathbf{A}\mathbf{y})$$

Optimization pipeline for spectral unmixing

Estimation error

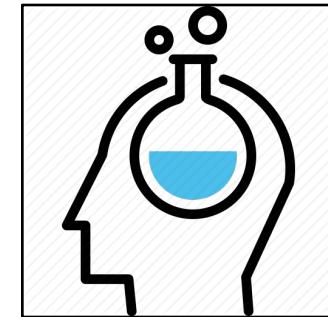
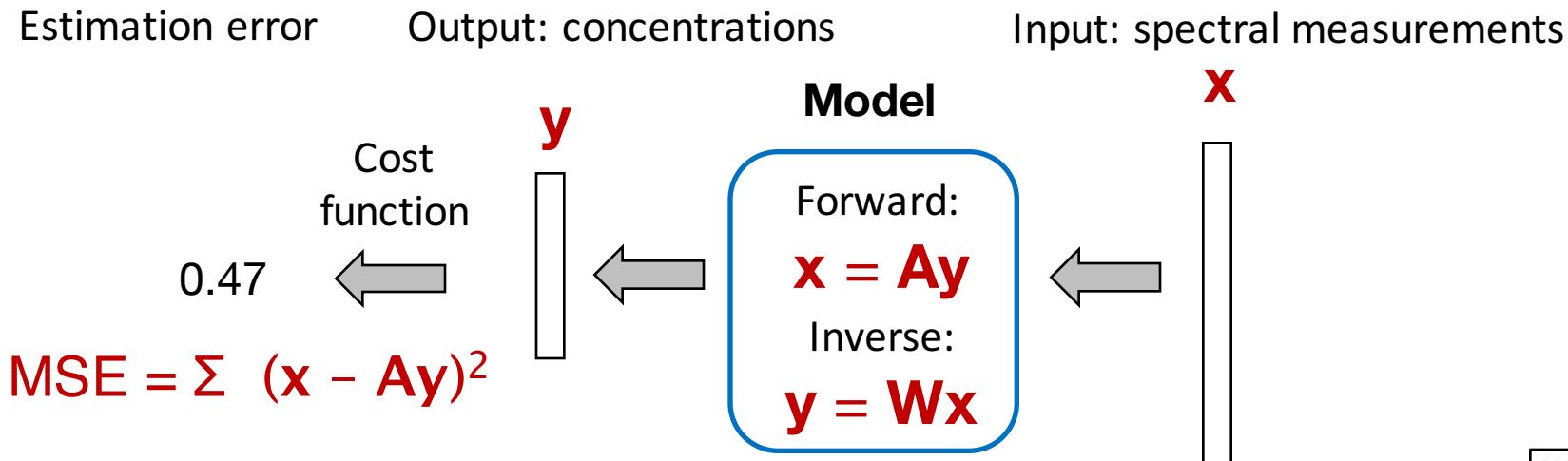
Output: concentrations

Input: spectral measurements



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

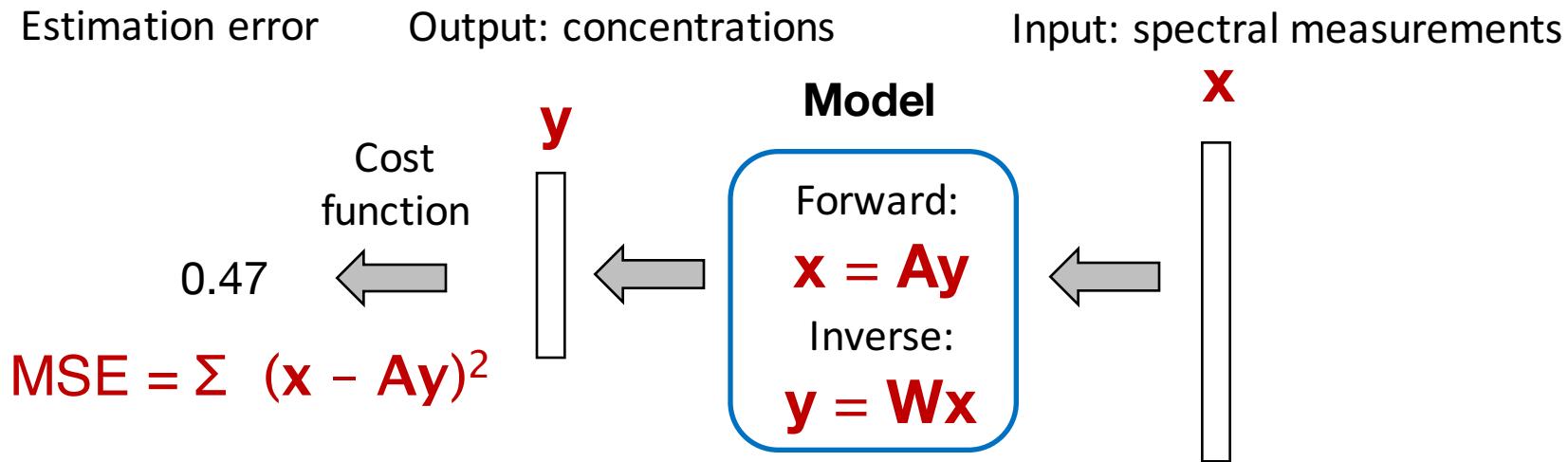
Optimization pipeline for spectral unmixing



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

A and W from first principles

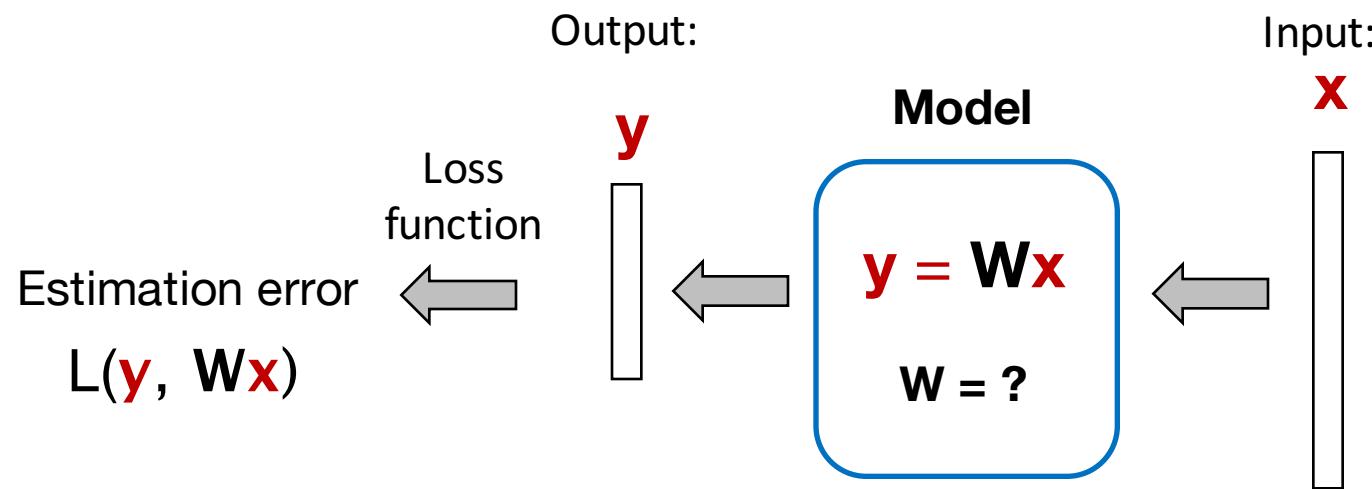
Optimization pipeline for spectral unmixing



Optimization: You only care about finding the best solution \mathbf{y}^*

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

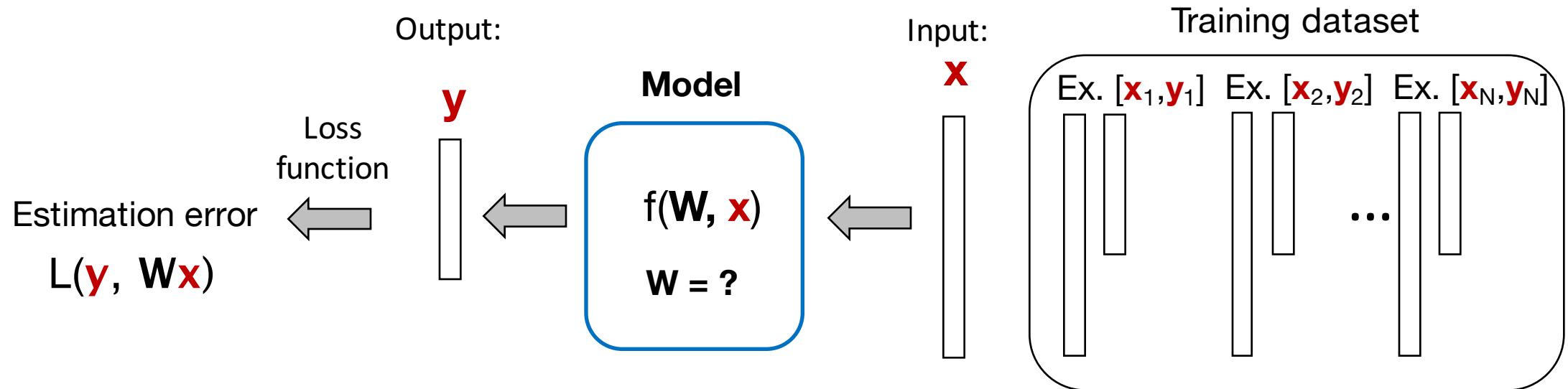
Pipeline for machine learning



Changes for machine learning framework:

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

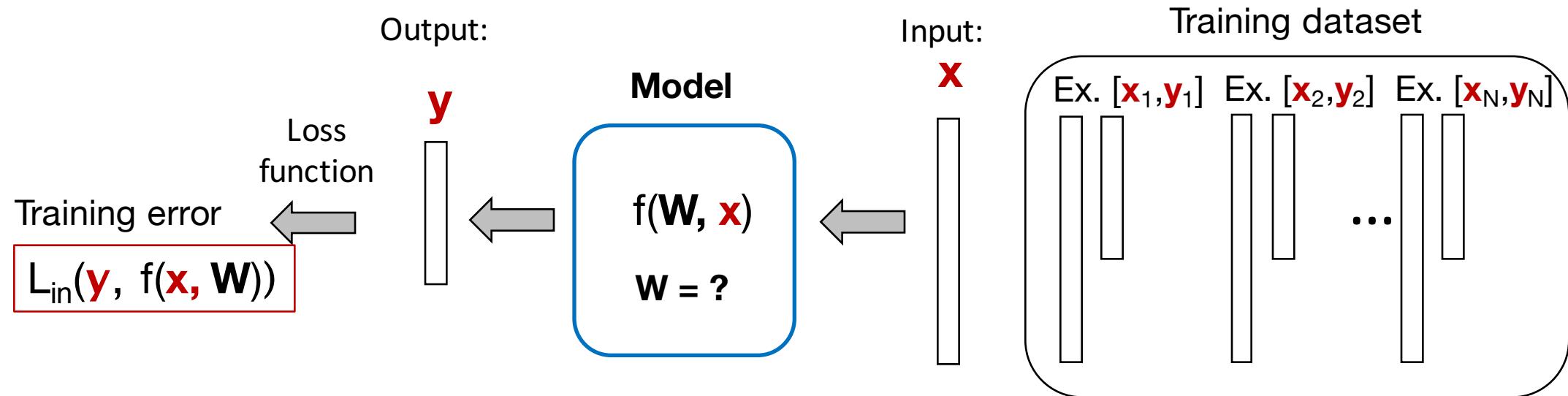
Pipeline for machine learning



Changes for machine learning framework:

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

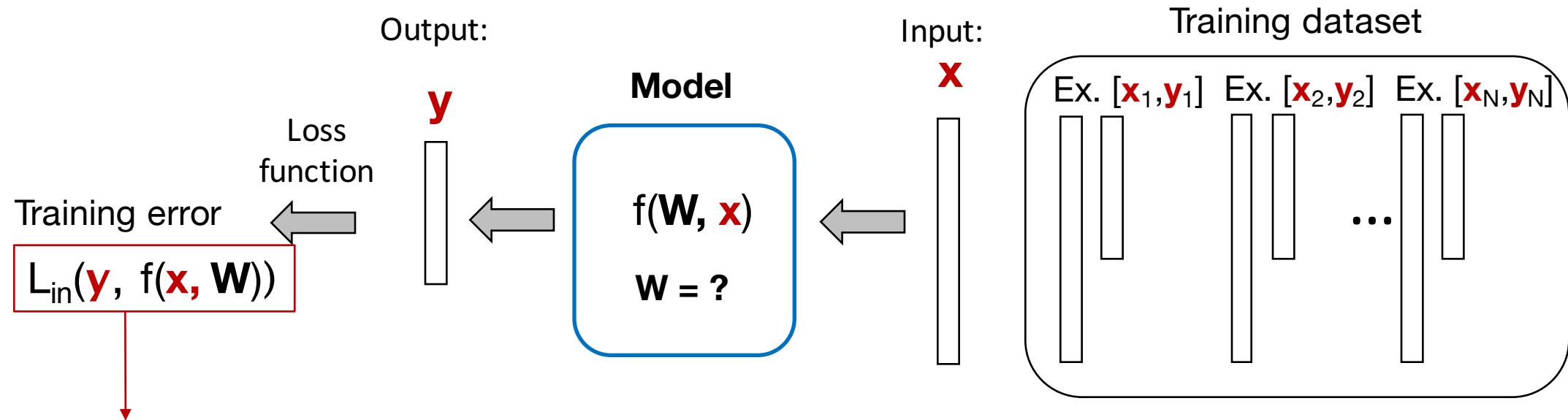
Pipeline for machine learning



Changes for machine learning framework:

1. Now just establish the mapping from inputs to outputs (here, matrix \mathbf{W})
2. Must first determine mapping $f(\mathbf{x}, \mathbf{W})$ using large set of “training” data
3. Use a *loss function* L that depends upon the training inputs (x, y) and the model (\mathbf{W})

Pipeline for machine learning

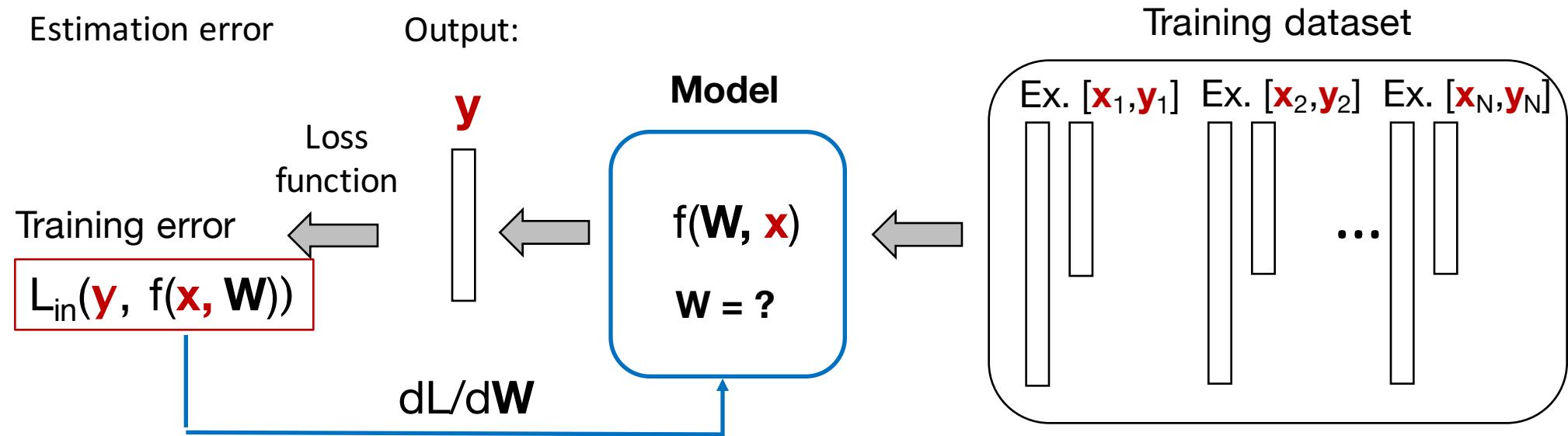


Training Error (“in class error”):

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

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

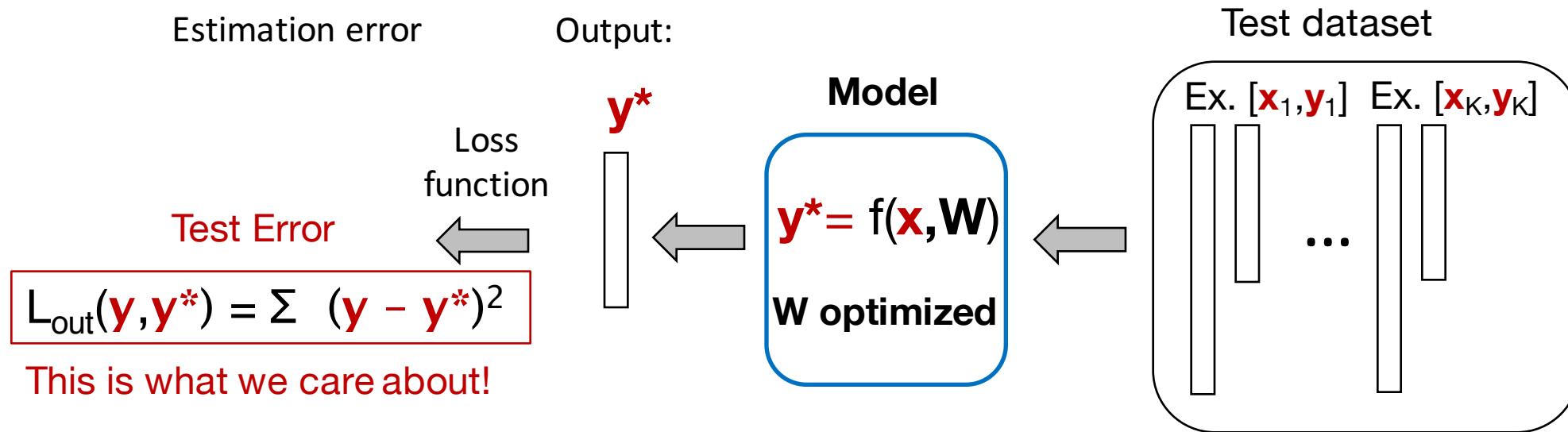
Pipeline for machine learning



Changes for machine learning framework:

1. Now just 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. 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 to guide gradient descent on L

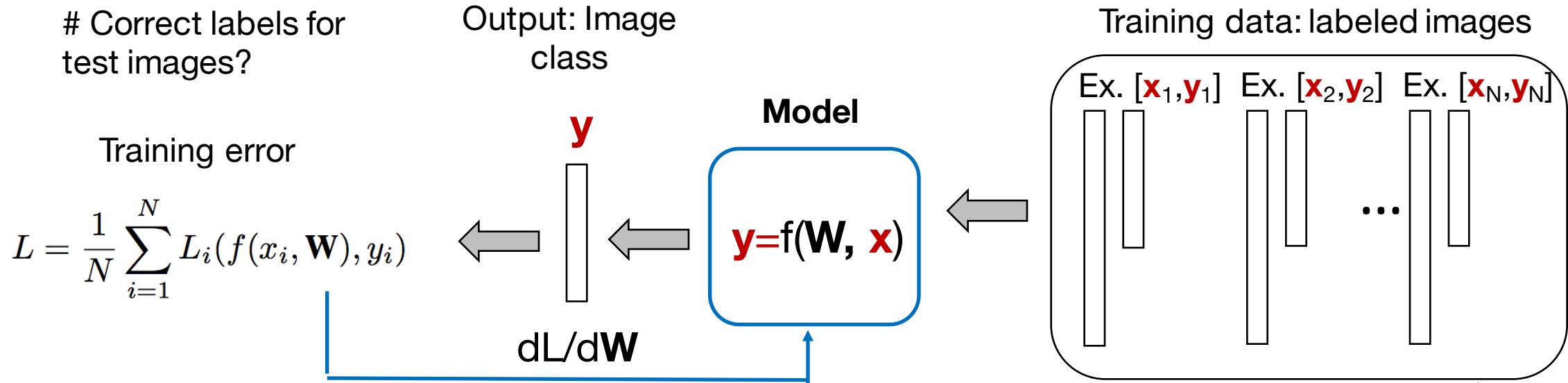
Pipeline for machine learning



Changes for machine learning framework:

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

Example: machine learning for image classification

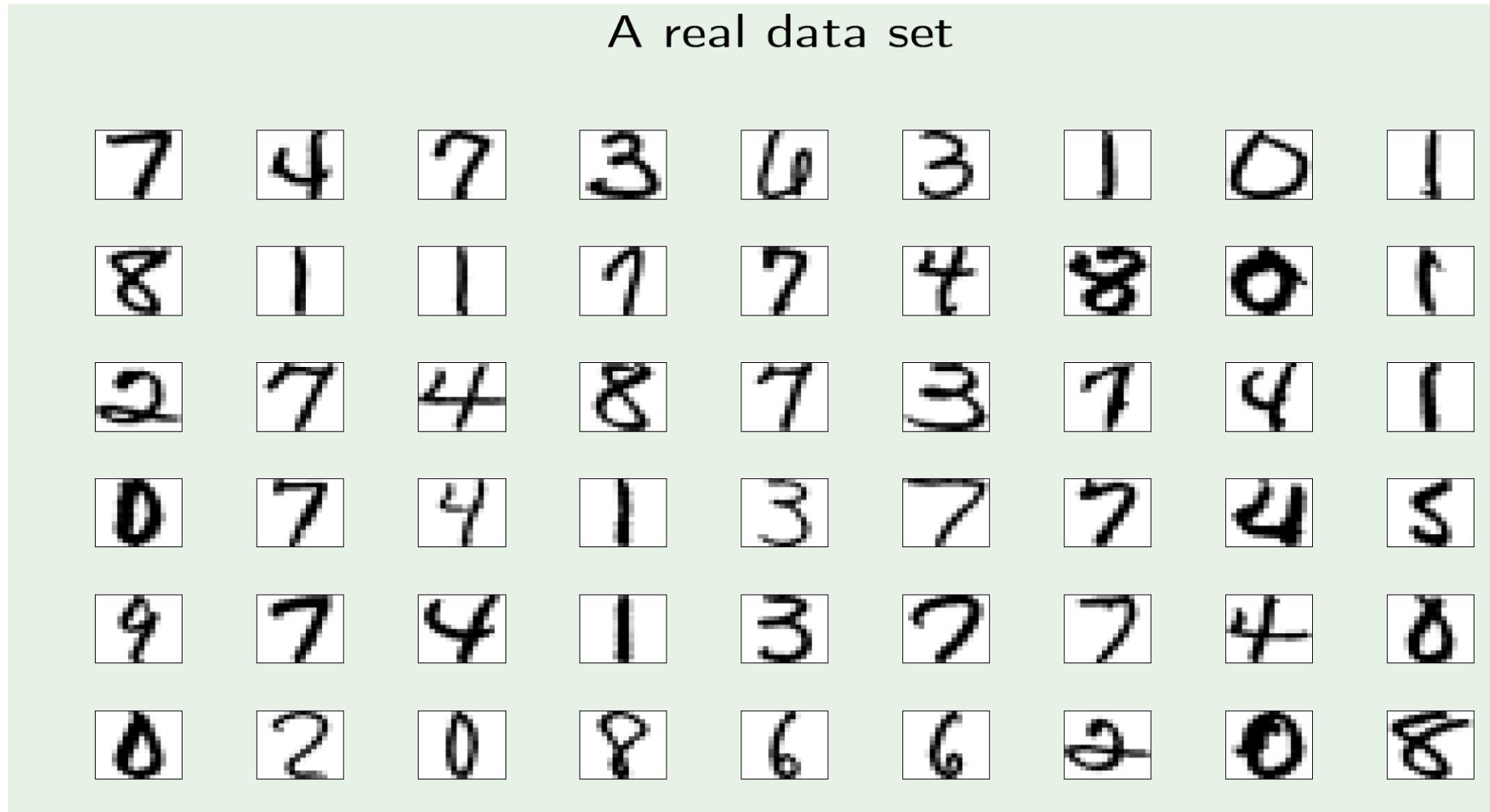


Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

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

Example: machine learning for image classification



Input representation

'raw' input $\mathbf{x} = (x_0, x_1, x_2, \dots, x_{256})$

linear model: $(w_0, w_1, w_2, \dots, w_{256})$

Features: Extract useful information, e.g.,

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

linear model: (w_0, w_1, w_2)

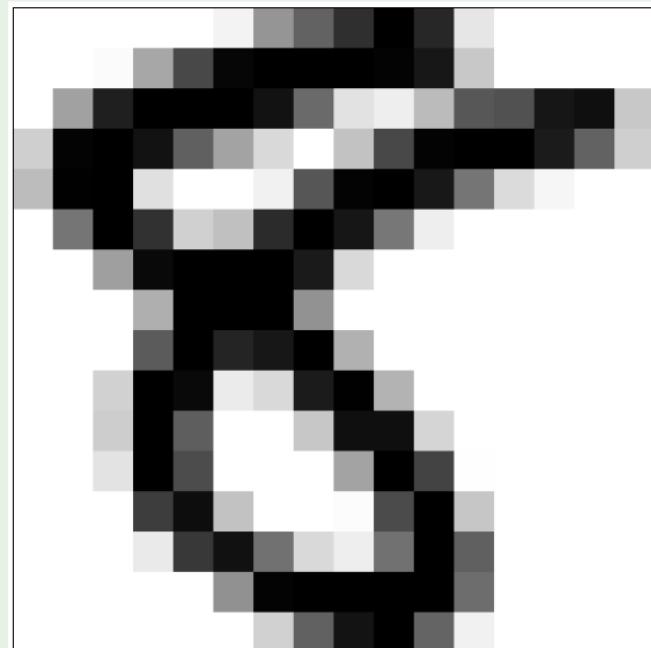
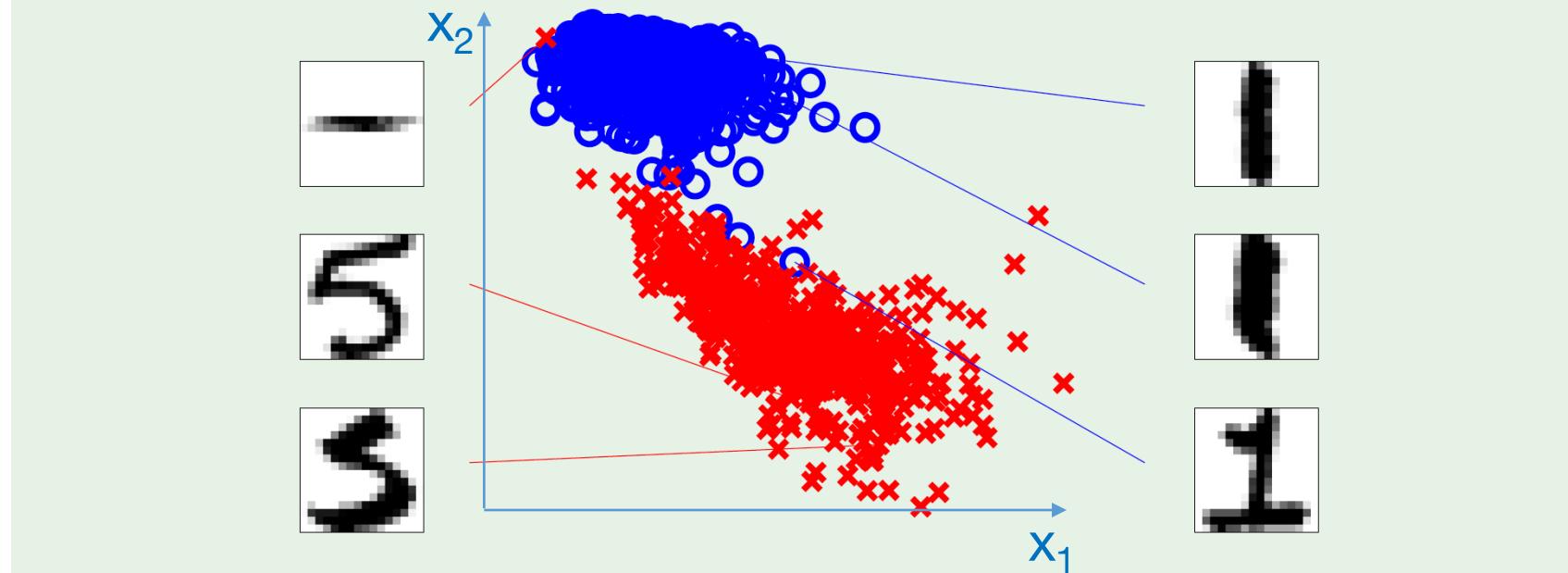


Illustration of features

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

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

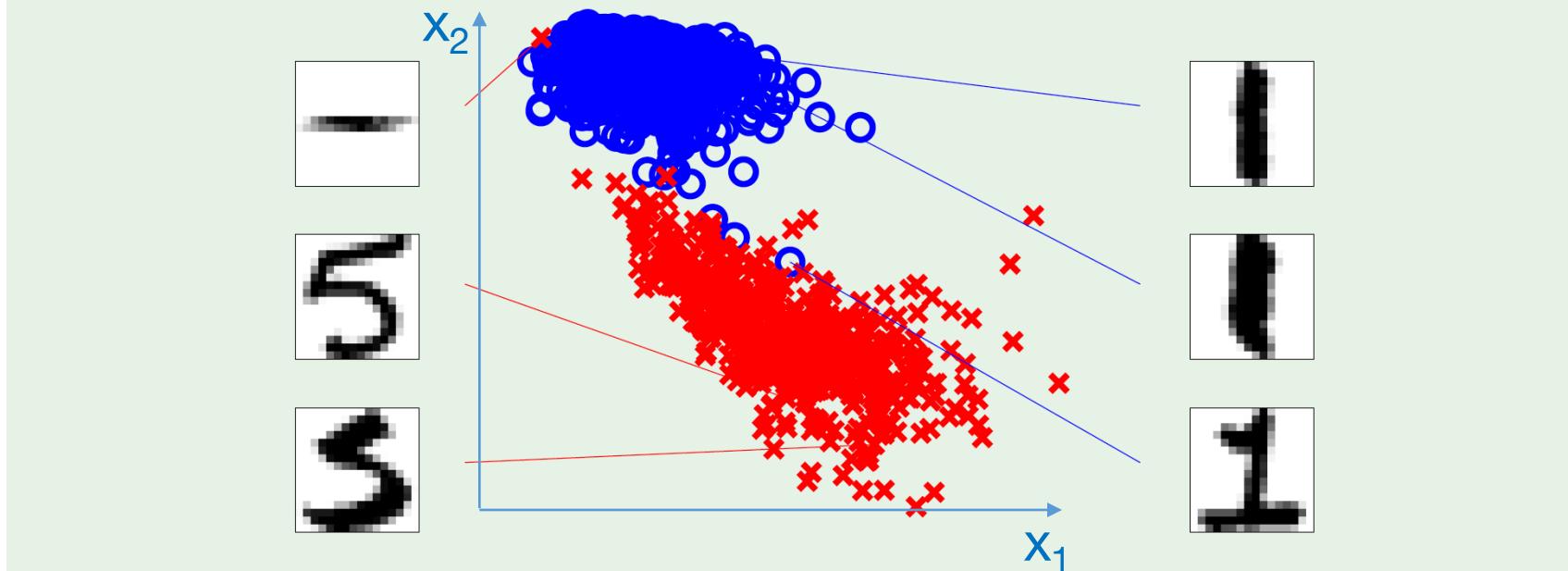
This transforms $wx+b$ into wx

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

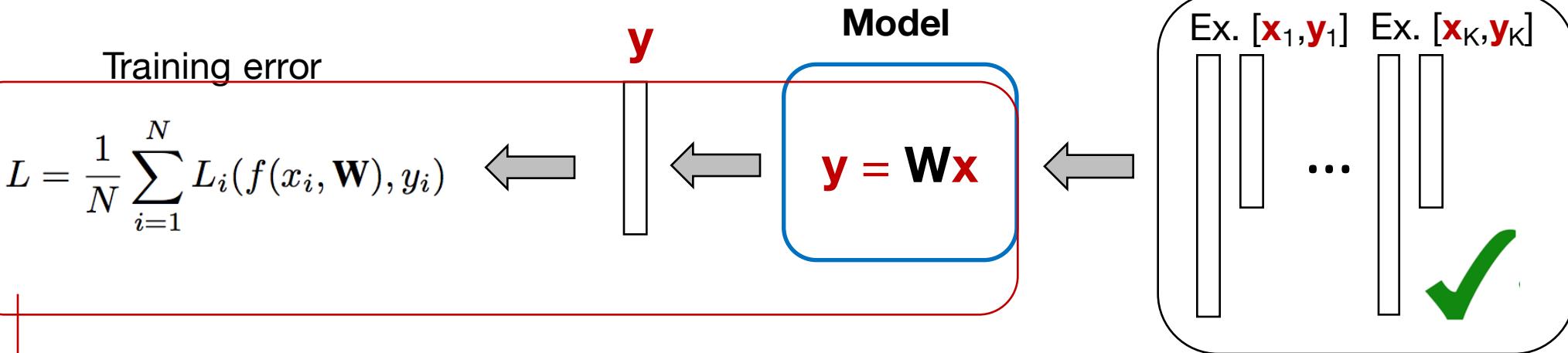
Test dataset

n=751-1000

Example: machine learning for image classification

Correct labels with
test dataset?

Output: Image
class



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

✓ 1. Labeled examples

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

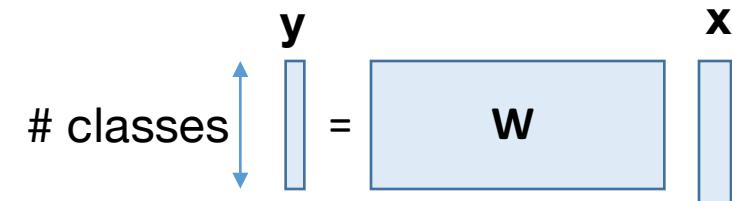
→ 2. A model and loss function

Let's start with a simpler approach: linear regression

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

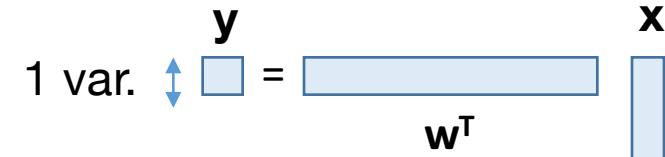
General linear model:

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



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

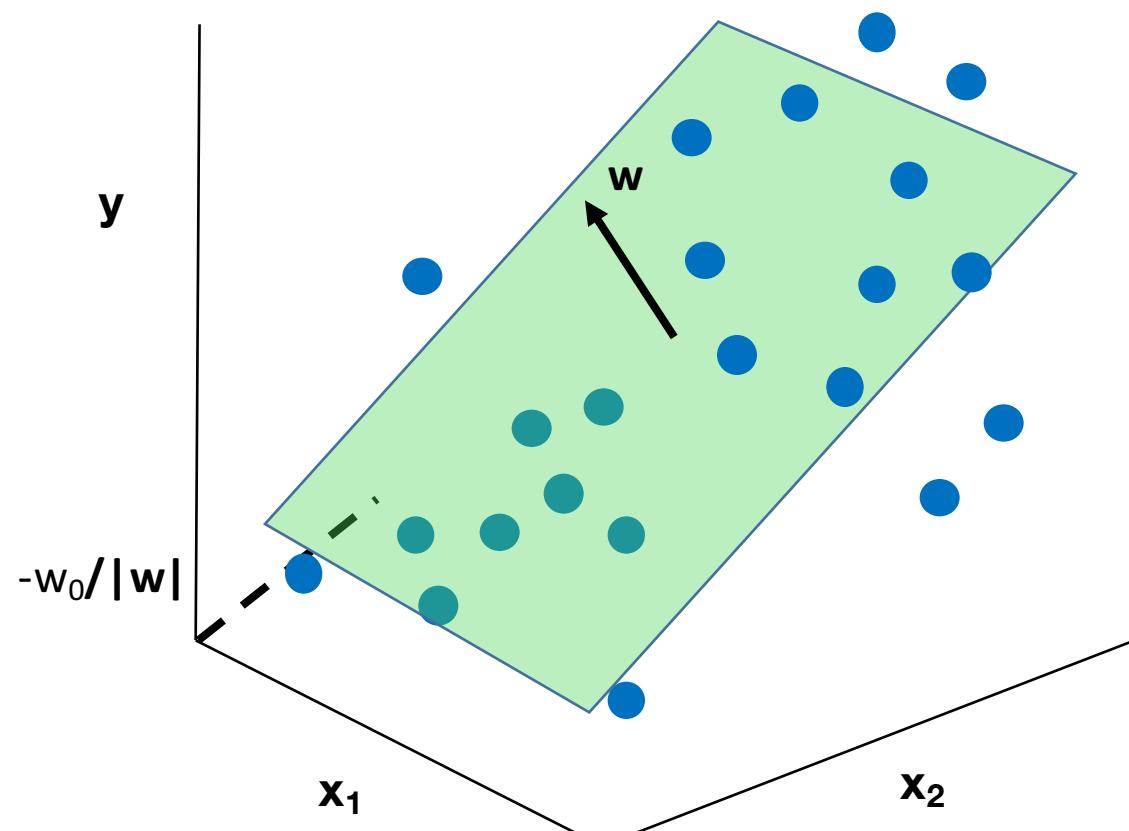
$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

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?



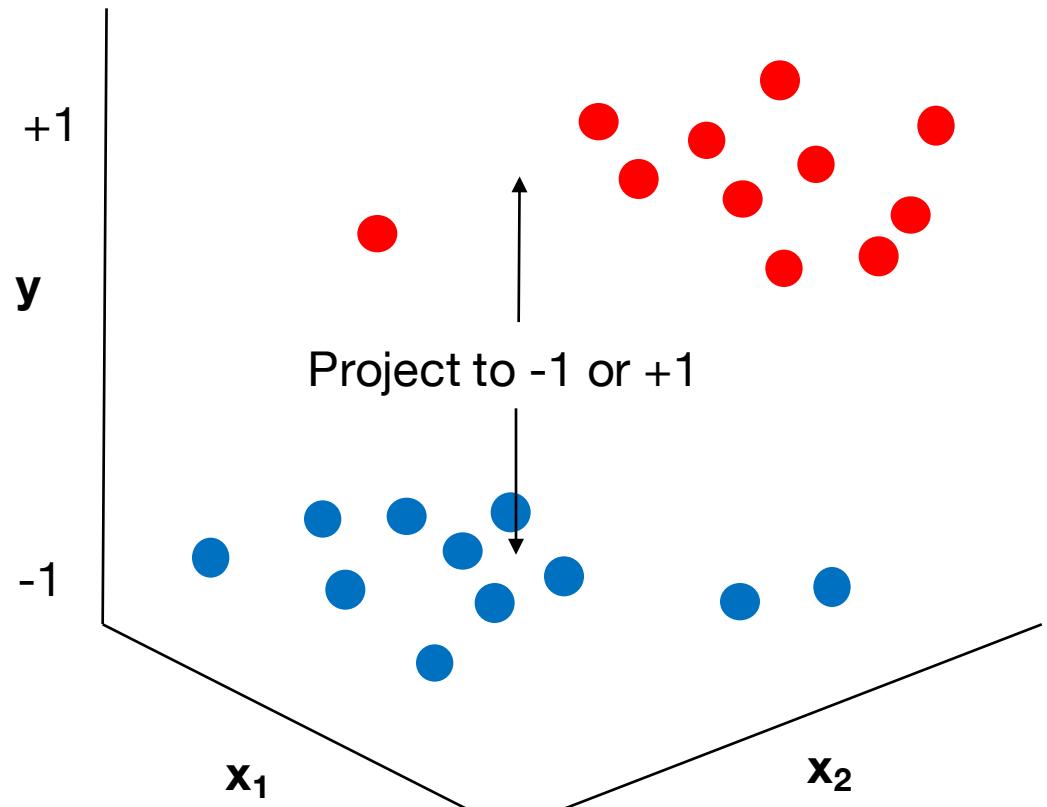
Without $\text{sgn}()$: regression for best fit

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

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

- 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?



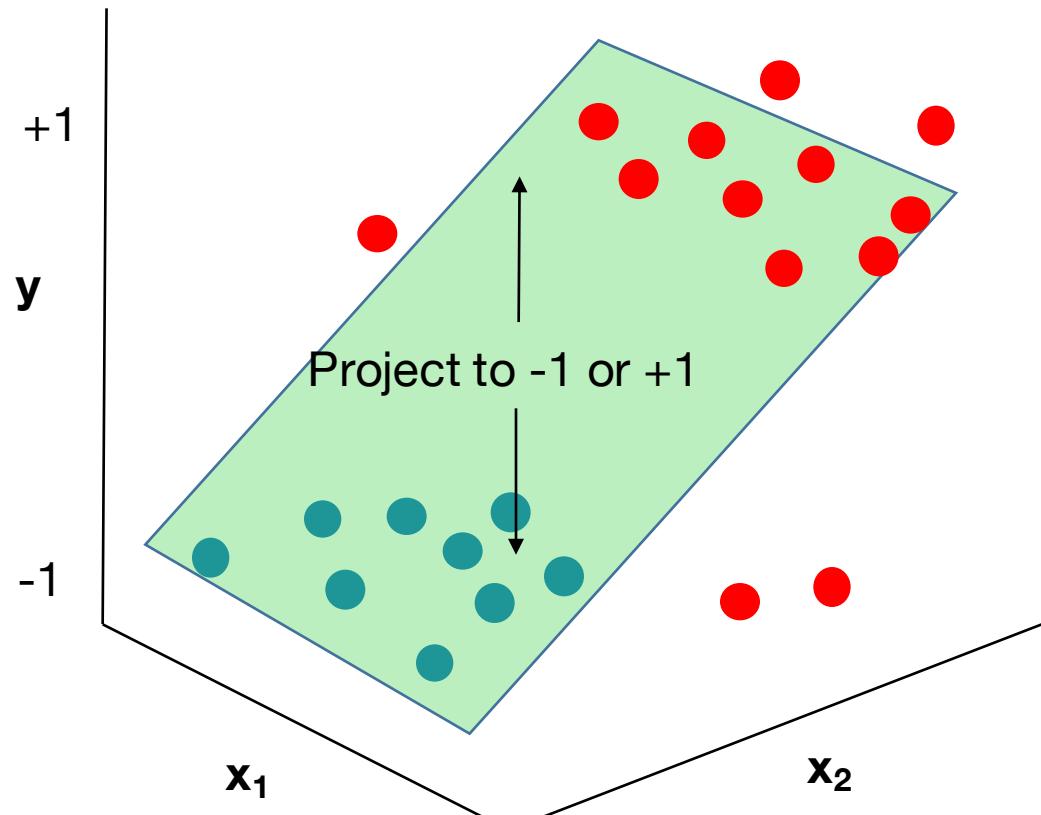
Without $\text{sgn}()$: regression for best fit

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

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

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

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



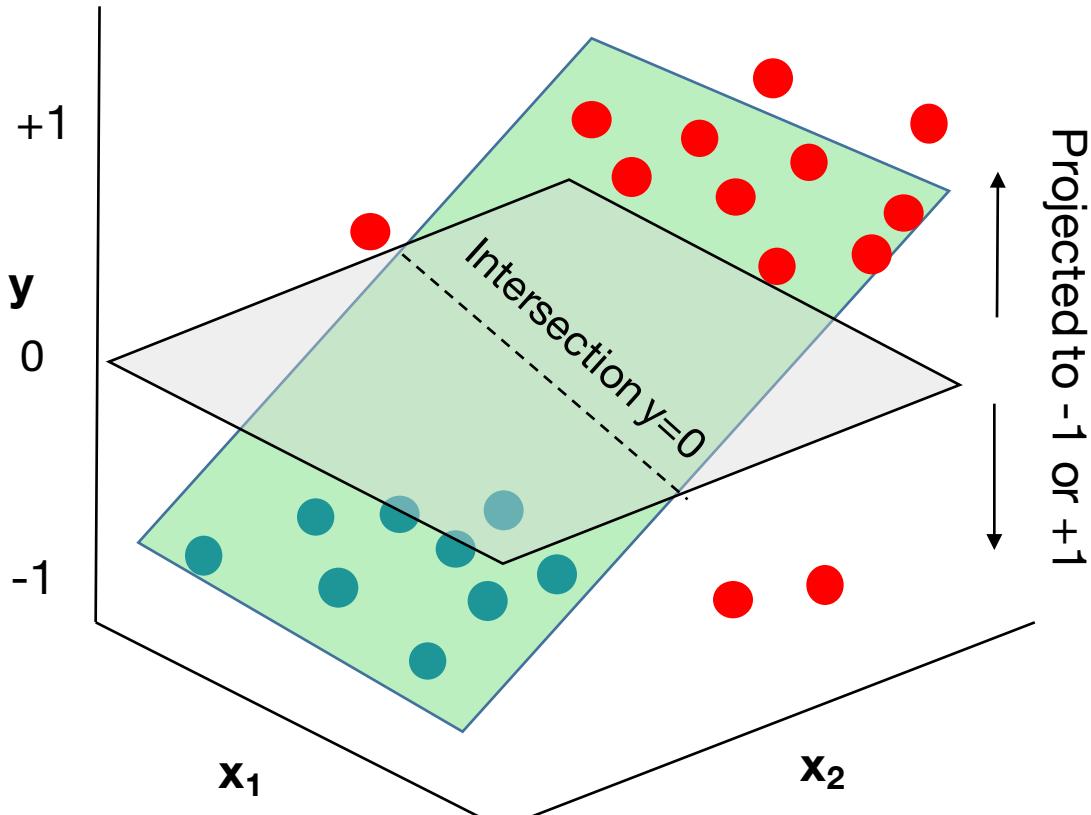
Without $\text{sgn}()$: regression for best fit

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

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

- 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?



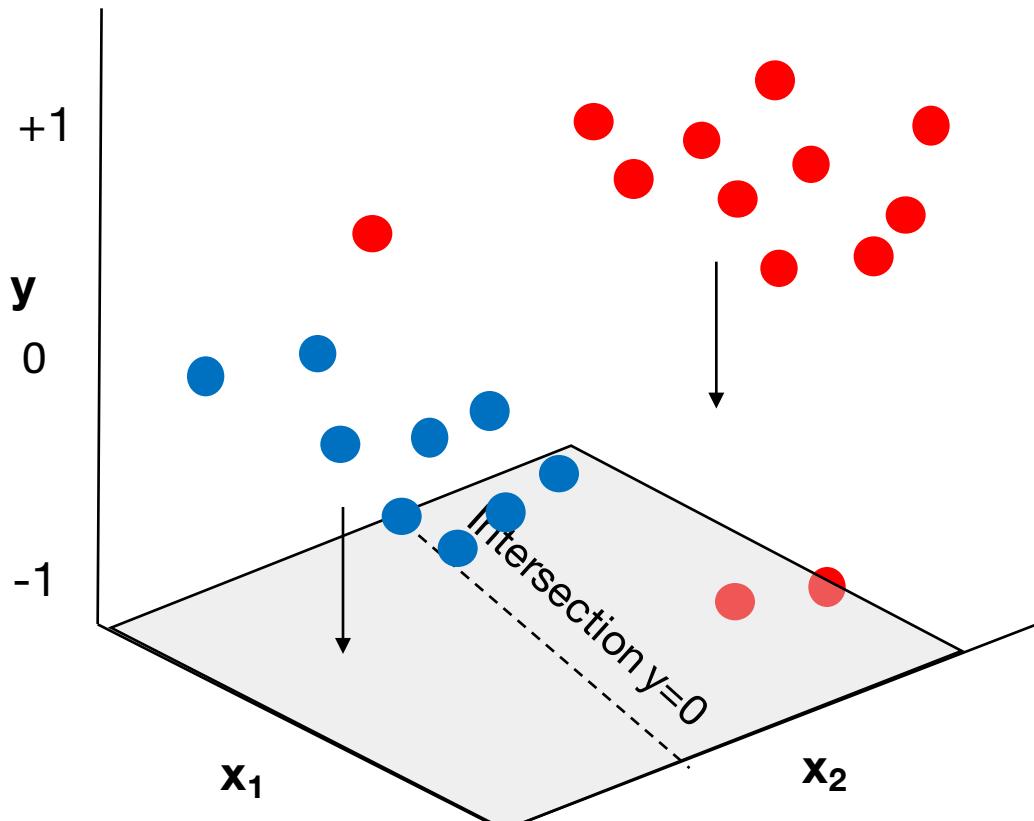
With $\text{sgn}()$ operation:

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

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

- 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?



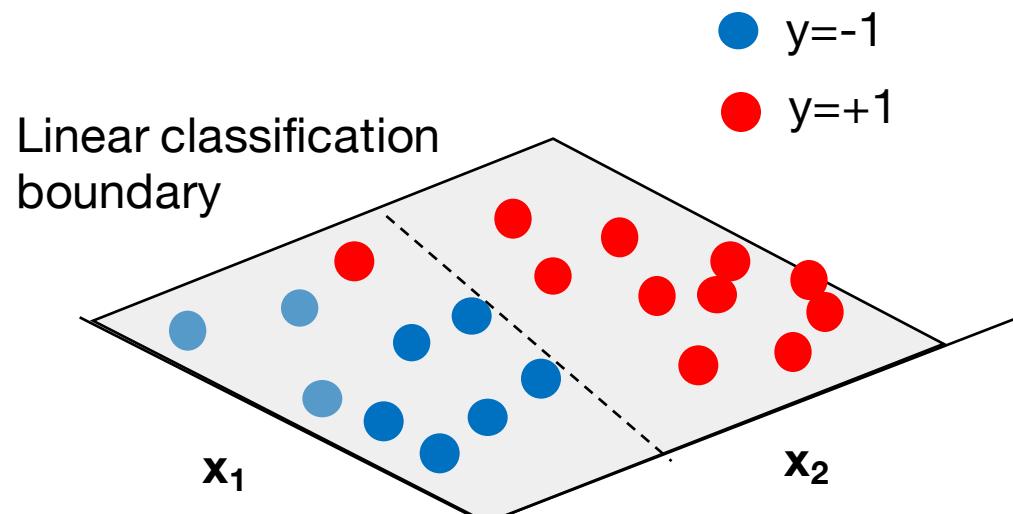
With $\text{sgn}()$ operation:

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

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

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

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



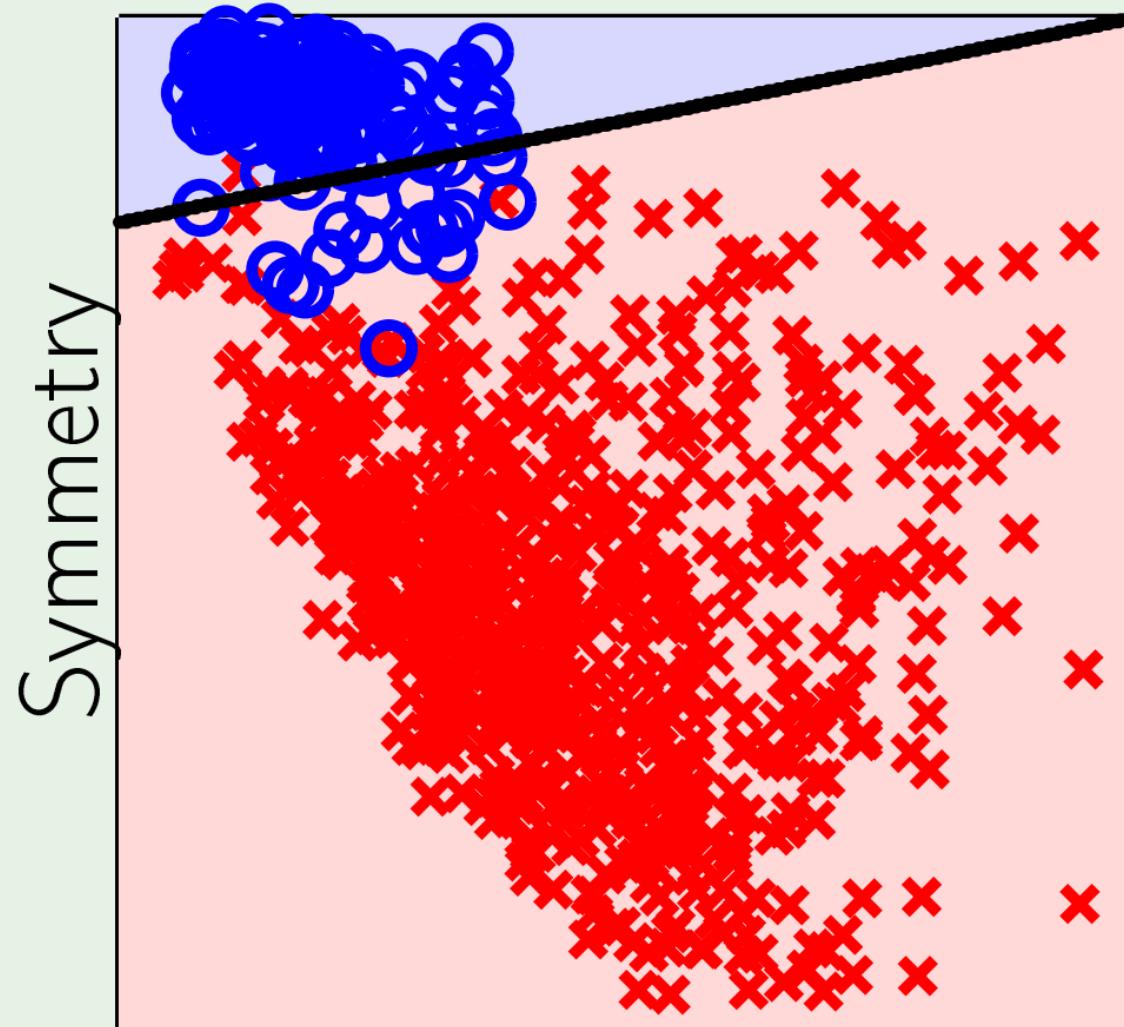
With $\text{sgn}()$ operation:

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

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

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

Linear regression boundary



Average Intensity

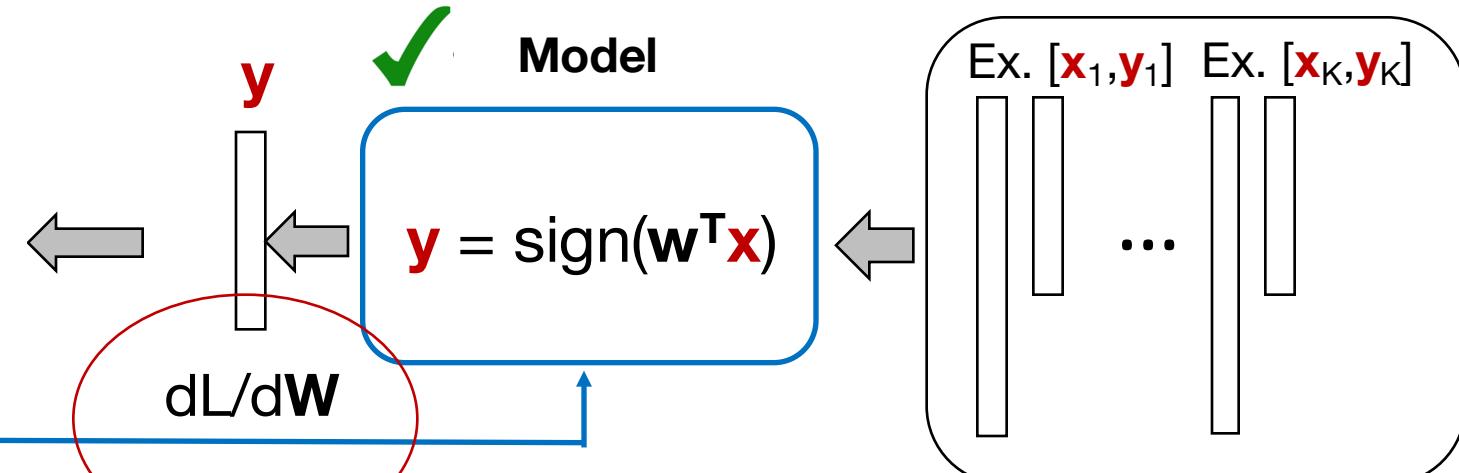
Example: machine learning for image classification

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)

$$L\left(\mathbf{w}\right) ~=~ \frac{1}{N}\sum_{n=1}^N\left(\mathbf{w}^\top\!\mathbf{x}_n - y_n\right)^2$$

$$\begin{aligned}
 L(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 \\
 &= \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2
 \end{aligned}$$

where

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Write this out as a matrix equation:

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

Note: Training data goes
into “dictionary” matrix

Take derivative wrt w and set to 0:

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

Solution is pseudo-inverse:

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

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

The linear regression algorithm

- 1: Construct the matrix \mathbf{X} and the vector \mathbf{y} from the data set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as follows

$$\mathbf{X} = \underbrace{\begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}}_{\text{input data matrix}}, \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

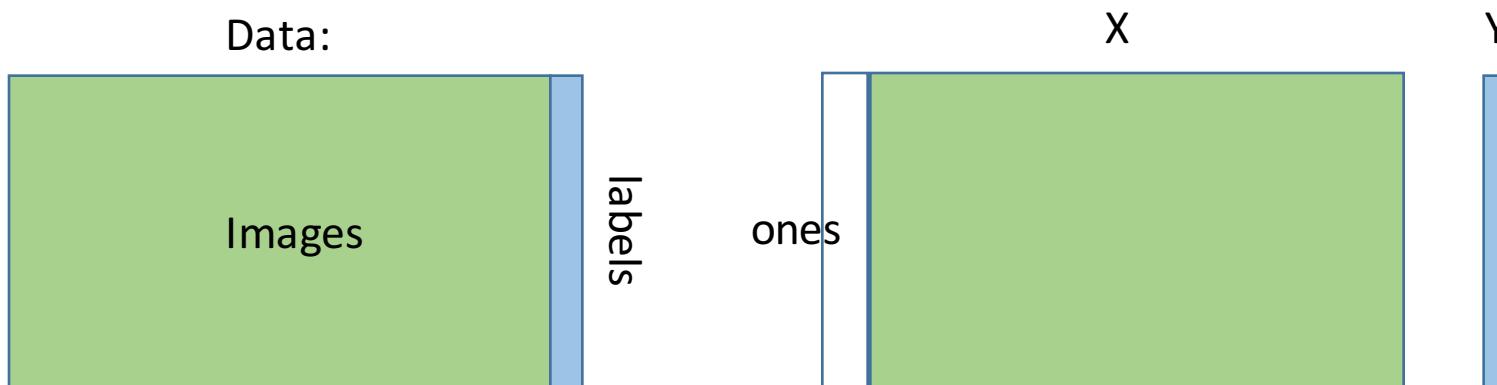
Download

- 2: Compute the pseudo-inverse $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

- 3: Return $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$.

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

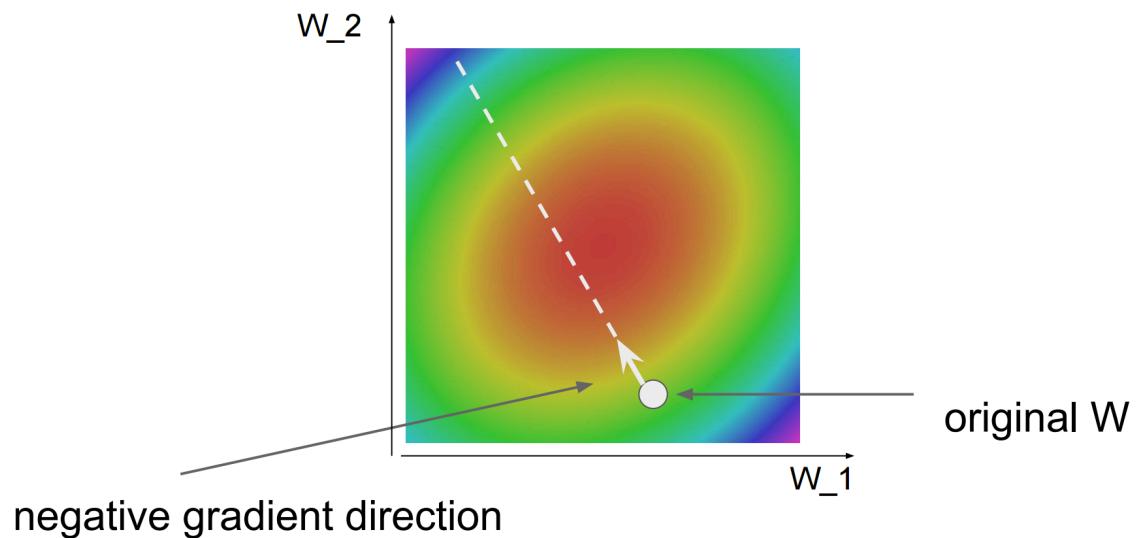
(harder)

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:

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

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)

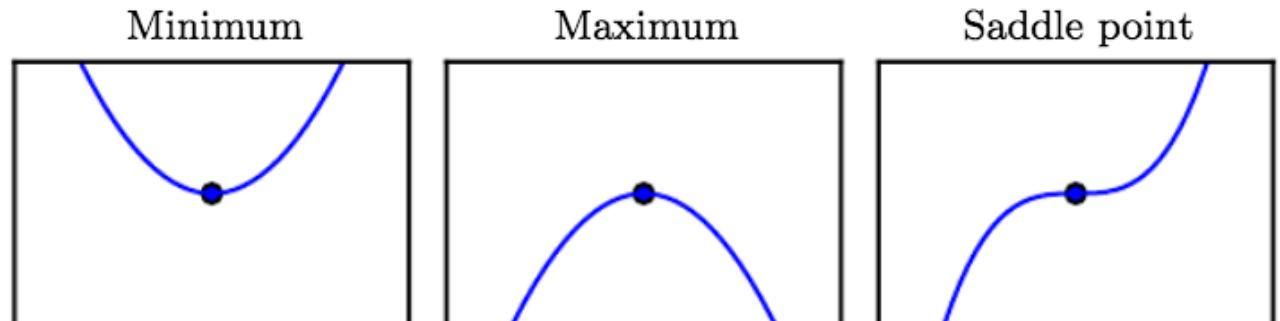
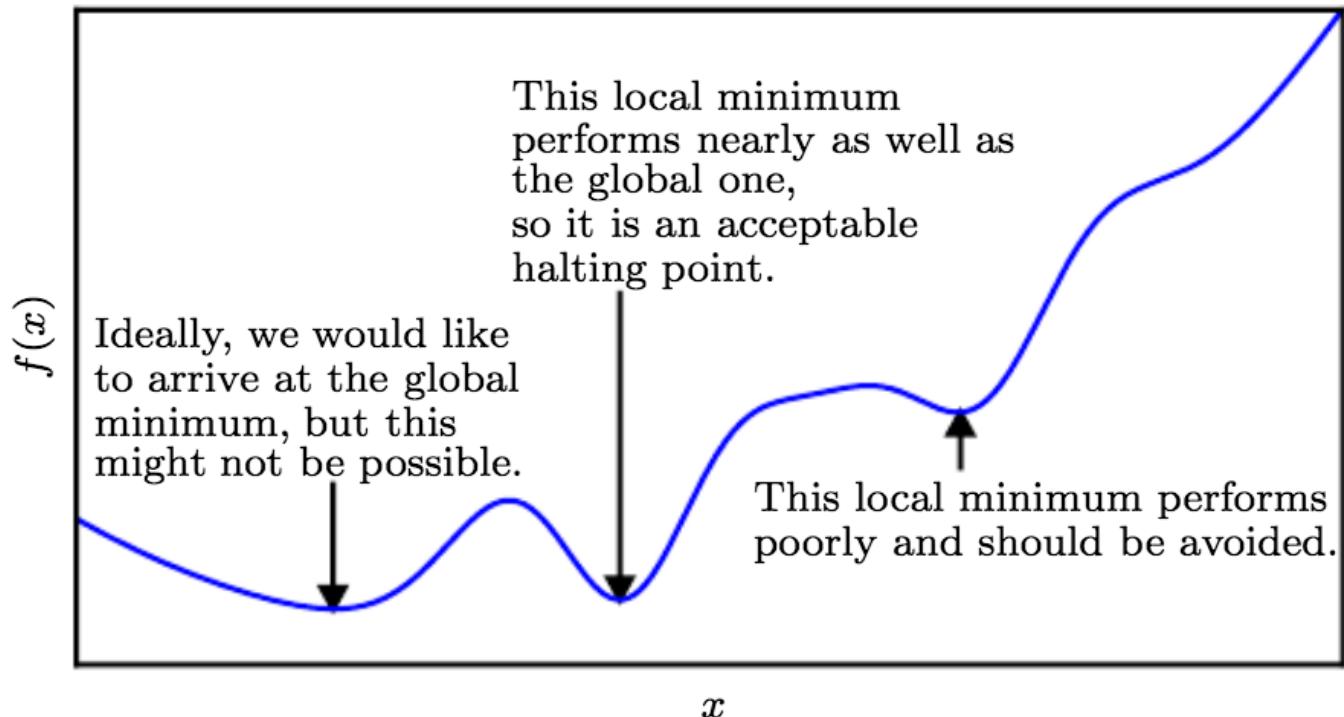
2. Numerical gradient descent

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

(harder)

Some quick details about gradient descent

- For non-convex functions, local minima can obscure the search for global minima
- Analyzing critical points (plateaus) of function of interest is important
- Critical points at $df/dx = 0$
- 2nd derivative d^2f/dx^2 tells us the type of critical point:
 - Minima at $d^2f/dx^2 > 0$
 - Maxima at $d^2f/dx^2 < 0$



Some quick details about gradient descent

Often we'll have functions of m variables

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{e.g., } f(\mathbf{x}) = \sum (\mathbf{Ax} - \mathbf{y})^2)$$

Will take partial derivatives $\frac{\partial}{\partial x_i} f(\mathbf{x})$ and put them in **gradient vector g**: $\nabla_{\mathbf{x}} f(\mathbf{x})$

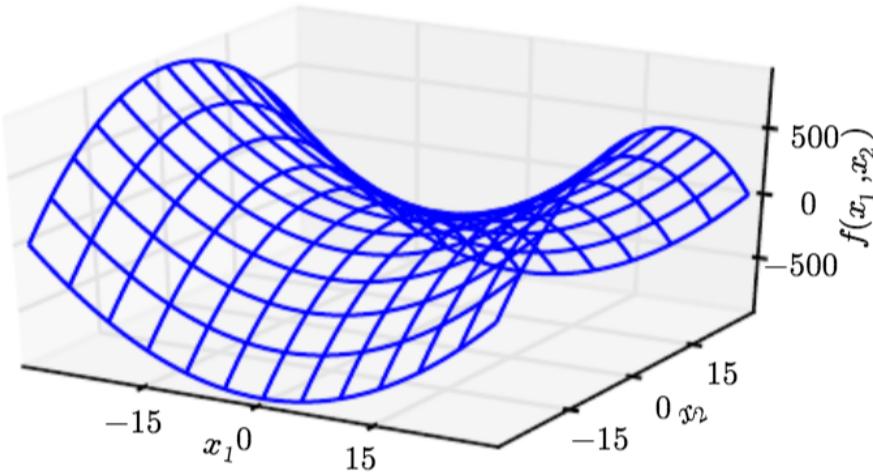
Will have many second derivatives: $\mathbf{H}(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$ **Hessian Matrix**

In general, we'll have functions that map m variables to n variables

$$\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n \quad (\text{e.g., } \mathbf{f}(\mathbf{x}) = \mathbf{Wx}, \mathbf{W} \text{ is } n \times m)$$

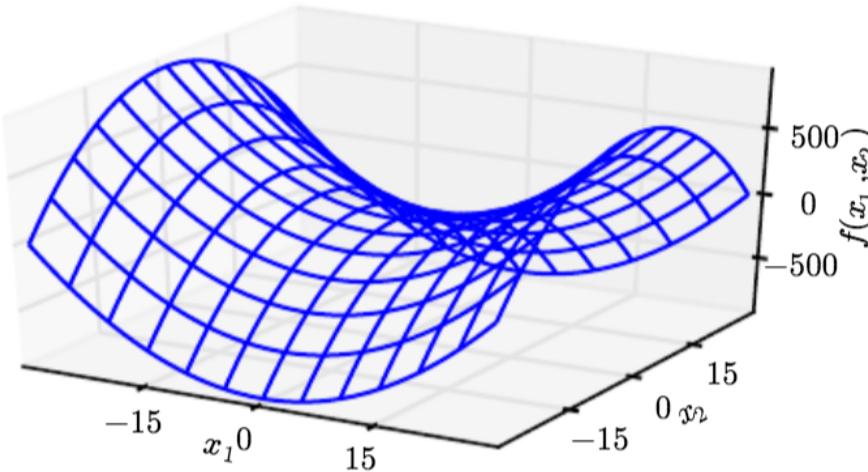
$$\mathbf{J} \in \mathbb{R}^{n \times m} \text{ of } \mathbf{f} : J_{i,j} = \frac{\partial}{\partial x_j} f(\mathbf{x})_i \quad \text{Jacobian Matrix}$$

Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

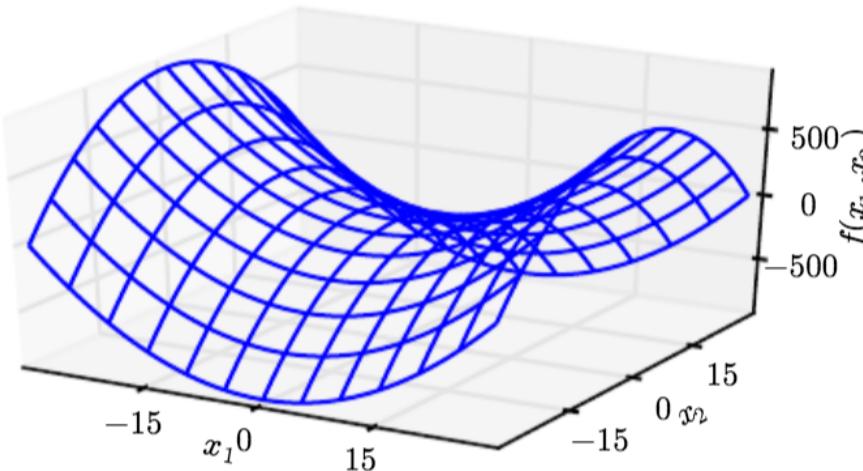
Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

$$\mathbf{g} = \begin{bmatrix} 2x_1 \\ -2x_2 \end{bmatrix}$$

Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

$$\mathbf{g} = \begin{bmatrix} 2x_1 \\ -2x_2 \end{bmatrix} \quad \mathbf{H} = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

- Convex functions have positive semi-definite Hessians (Trace ≥ 0)
- Trace/eigenvalues of Hessian are useful evaluate critical points & guide optimization

Steepest descent and the best step size ϵ

1. Evaluate function $f(\mathbf{x}^{(0)})$ at an initial guess point, $\mathbf{x}^{(0)}$
2. Compute gradient $\mathbf{g}^{(0)} = \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$
3. Next point $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \epsilon^{(0)} \mathbf{g}^{(0)}$
4. Repeat – $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \epsilon^{(n)} \mathbf{g}^{(n)}$, until $|\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}| < \text{threshold } t$

```
while previous_step_size > precision and iters < max_iters:  
    prev_x = cur_x  
    cur_x -= gamma * df(prev_x)  
    previous_step_size = abs(cur_x - prev_x)  
    iters+=1
```

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Then, evaluate at the next step:

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Then, evaluate at the next step:

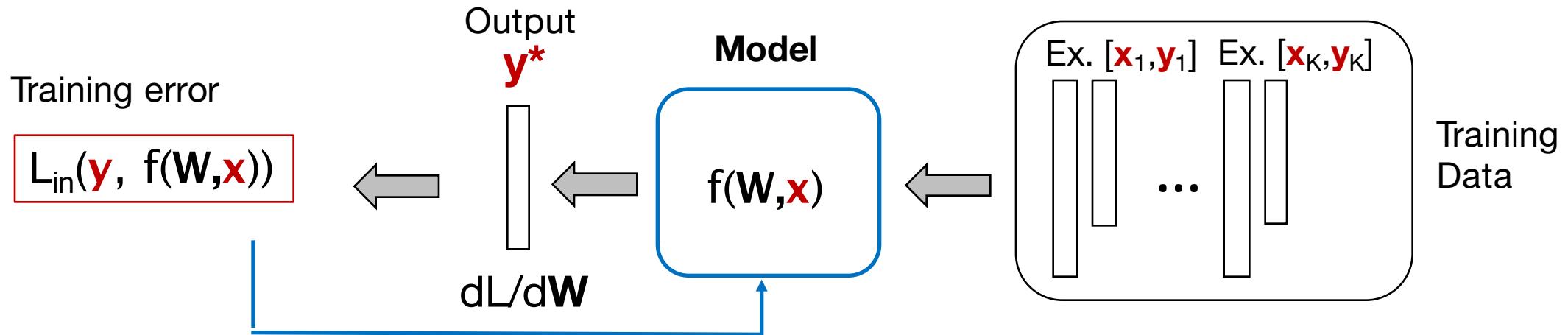
$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

Solve for optimal step (when Hessian is positive):

$$\boxed{\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}}.$$

Summary of machine learning pipeline:

1. Network Training

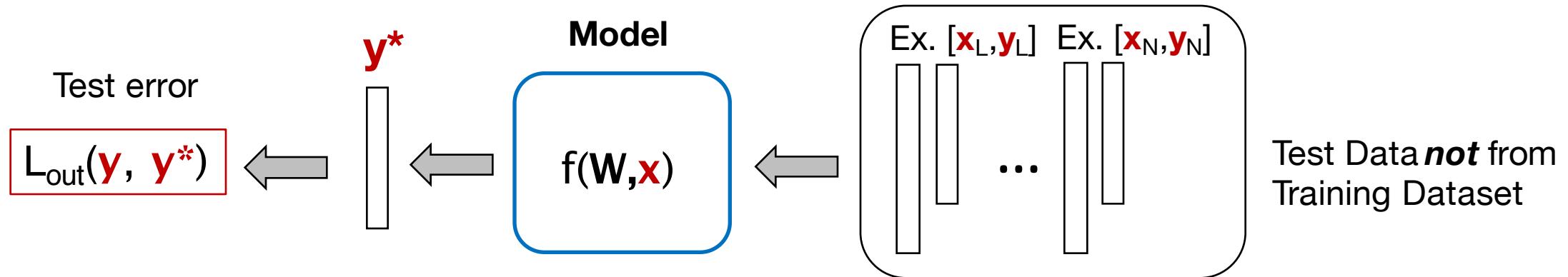


What we need for network training:

- 1. Labeled examples**
- 2. A model and loss function**
- 3. A way to minimize the loss function L**

Summary of machine learning pipeline:

2. Network Testing



What we need for network testing:

4. ***Unique* labeled test data**
5. **Evaluation of model error**