

A physical scientist's introduction to Machine Learning

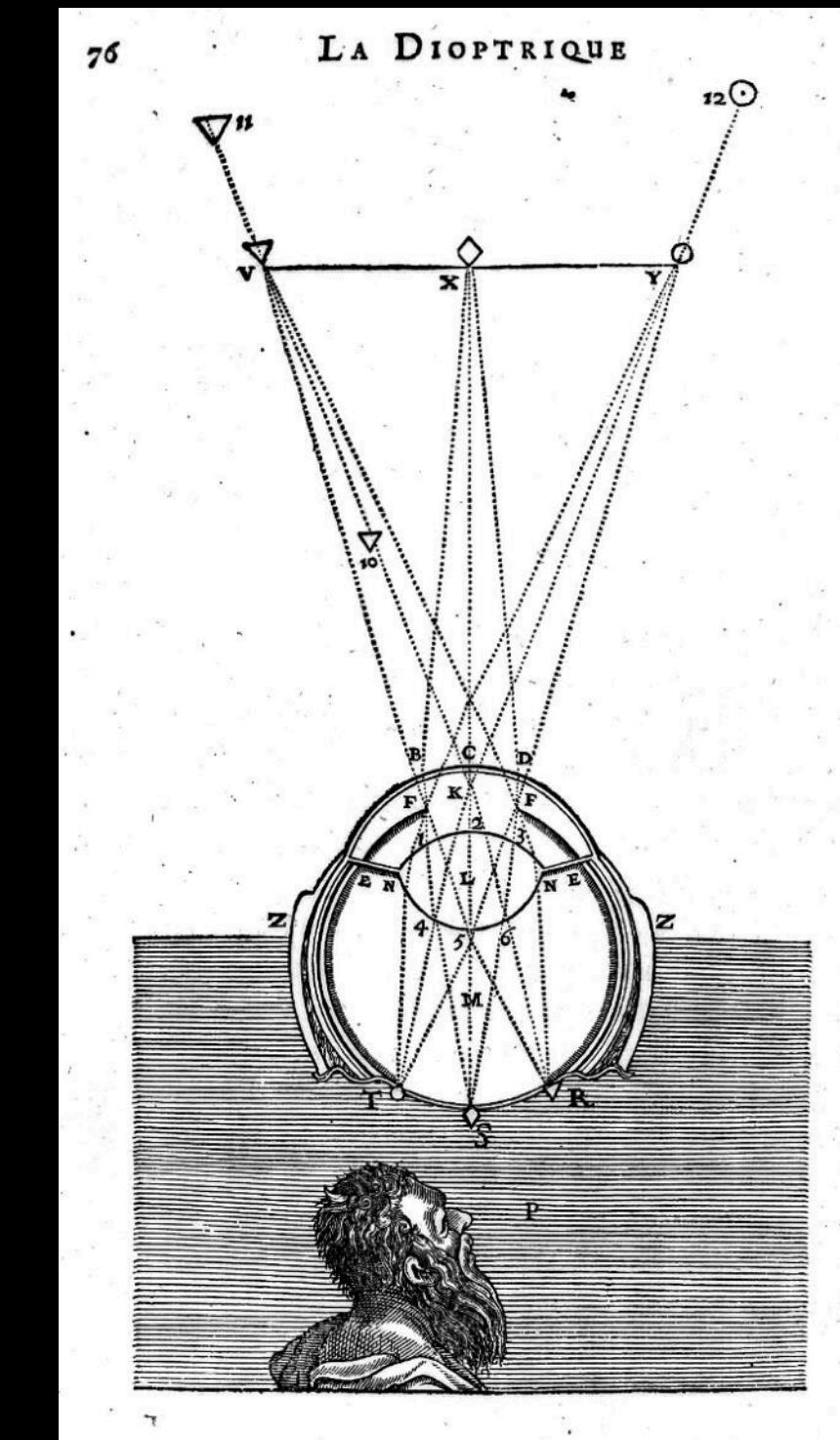
Physics from data



Galileo



Kepler



Descartes

Finding **constants** of nature that **generalize** in space and time

Laws are linear

Pascal's law (1653)



Hooke's law (1678)



Newton's law of viscosity (1701)



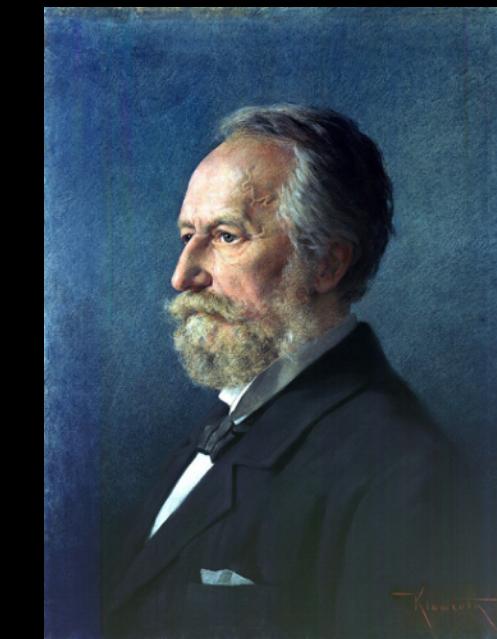
Ohm's law (1781)



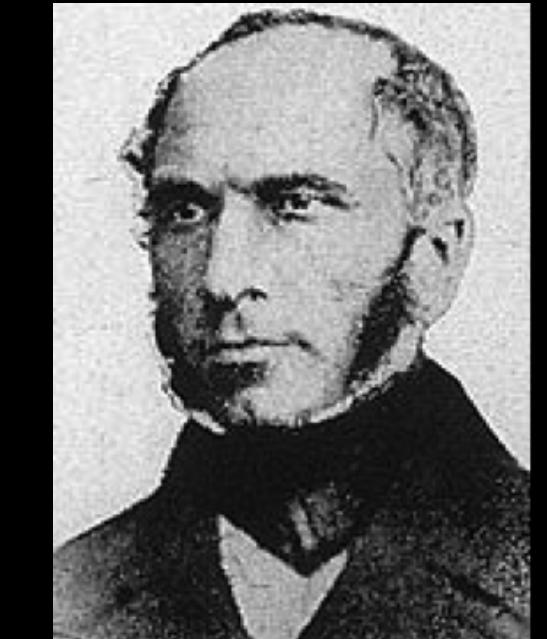
Fourier's law (1822)



Fick's law (1855)



Darcy's law (1856)



$$\Delta p = \rho g \Delta h$$

$$F = -kx$$

$$\tau = \mu \frac{du}{dy}$$

$$I = V/R$$

$$q = -k \frac{dT}{dx}$$

$$J = -D \frac{dC}{dx}$$

$$Q = \frac{kA}{\mu L} \Delta p$$

Ideal gas law (1834)

Amonton's law (1808)

Charles's law (1787)

Boyle's law (1662)

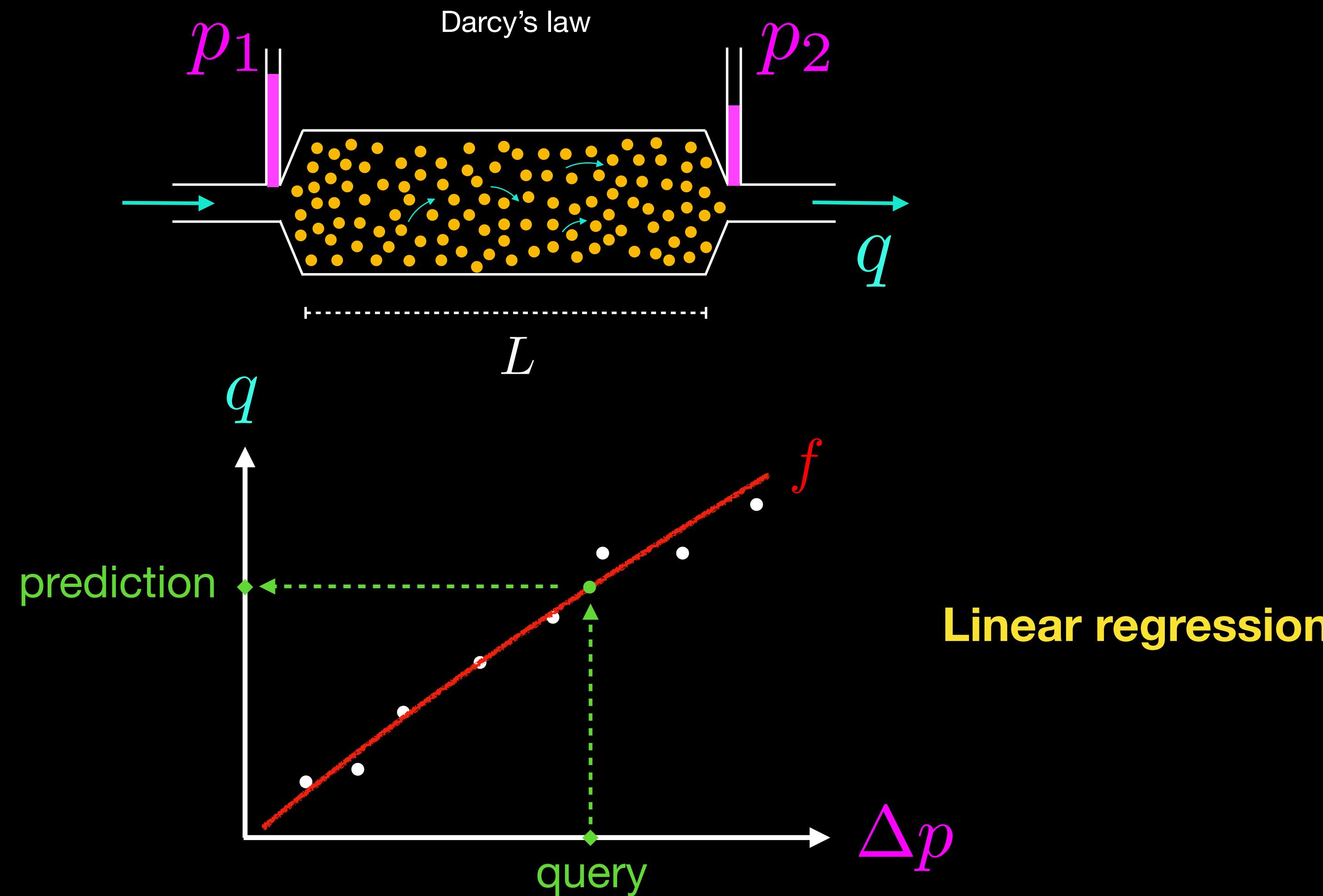
Avogadro's law (1811)

$$\frac{PV}{TN} = k_B$$

From experiment to Law

p_1	p_2	q
1.3	1.0	22
1.6	1.5	23
3.4	2.4	46
4.8	3.5	67
6.7	4.5	83
...		
2.3	1.4	?

$$(p_1, p_2) \rightarrow f \rightarrow q$$



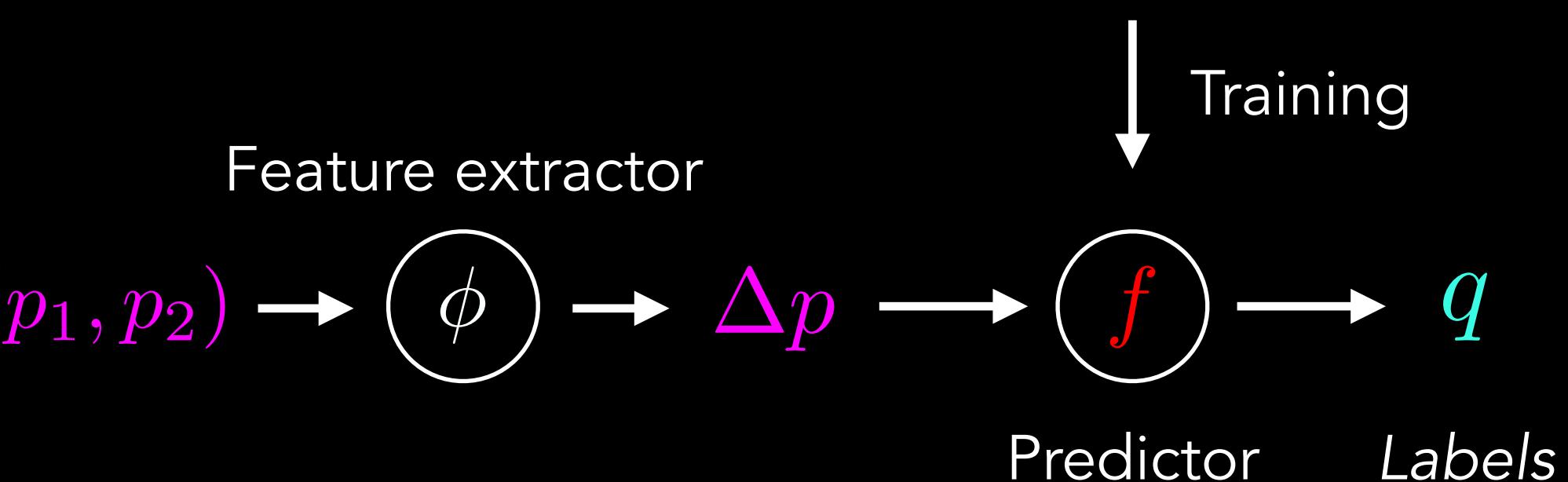
Machine or human learning

Training data: $\mathcal{D}_{\text{train}}$

Example →

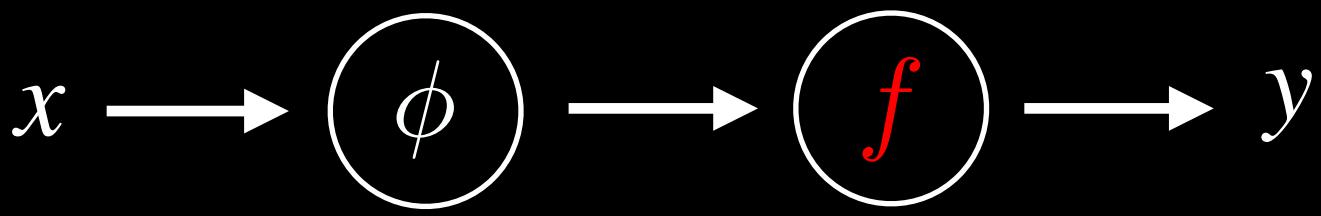
p_1	p_2	q
1.3	1.0	22
1.6	1.5	23
3.4	2.4	46
4.8	3.5	67
6.7	4.5	83
...		

- Which predictors are possible?
- How good is the predictor?
- How can we find the best predictor?



Linear predictors

(x_1, y_1)
(x_2, y_2)
(x_3, y_3)
(x_4, y_4)
\vdots
(x_n, y_n)



feature vector

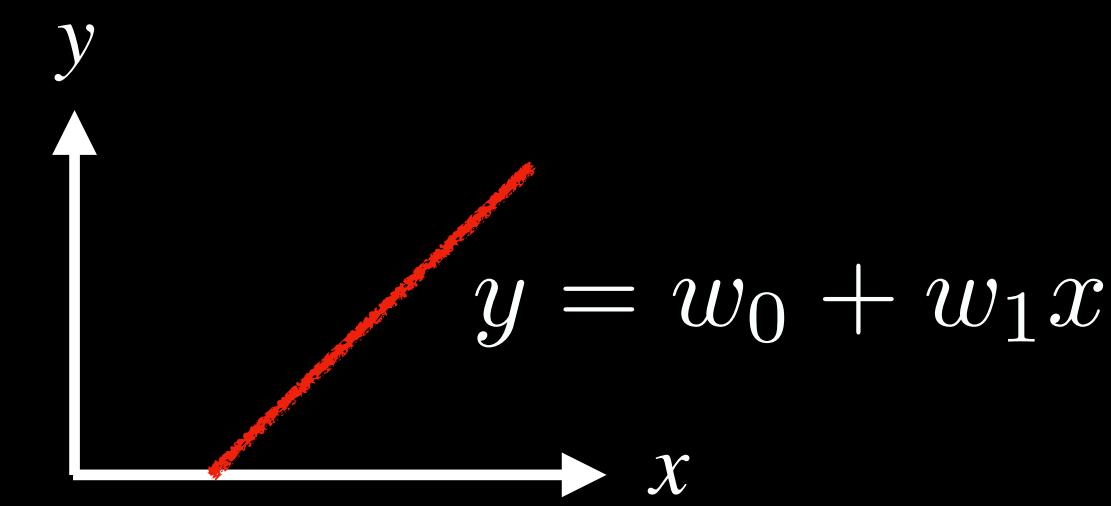
$$\phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_d(x)]$$

linear predictor

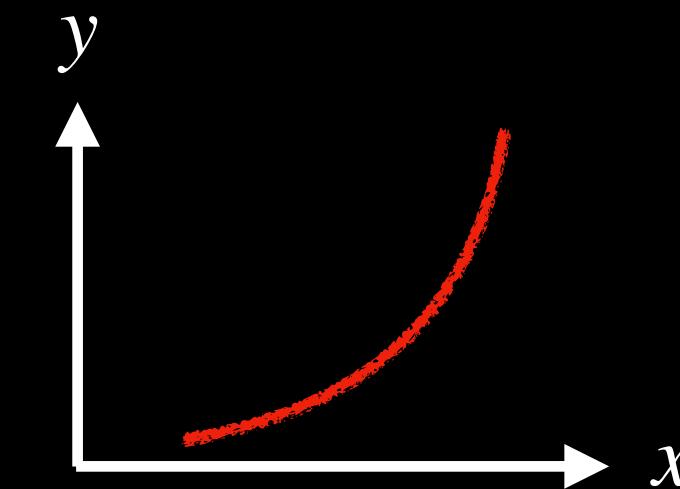
$$\begin{aligned}f_{\mathbf{w}} &= \mathbf{w} \cdot \phi(x) \\&= w_1\phi_1(x) + w_2\phi_2(x) + \dots + w_d\phi_d(x)\end{aligned}$$

What is linear?

Adding a bias: $\phi(x) = [1, x]$

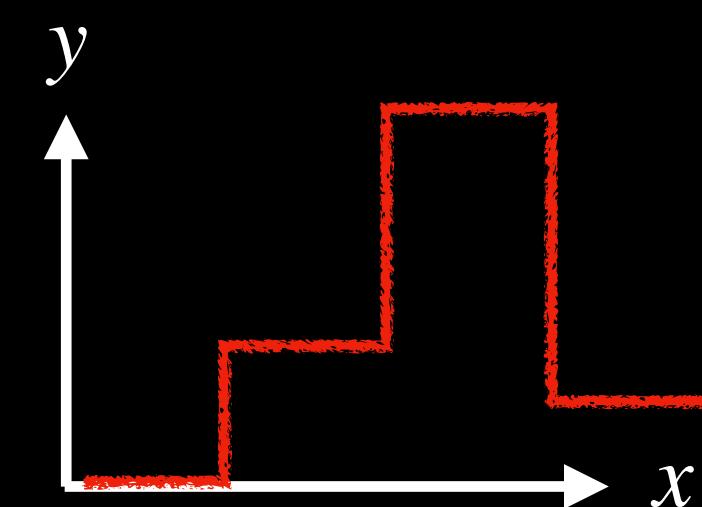


A polynomial predictor: $\phi(x) = [1, x, x^2, x^3]$



Sines and cosines: $\phi(x) = [1, x, \sin(3x)]$

Indicator functions: $\phi(x) = [\mathbf{1}[0 < x \leq 1], \mathbf{1}[1 < x \leq 2], \mathbf{1}[2 < x \leq 3]]$



Derivatives: $\phi(x, t) = \left[x^2, \frac{\Delta x}{\Delta t} \right]$

Linear in $\phi(x)$ and w but not in x

Hypothesis Class

$$\mathcal{F} = \{ f_{\mathbf{w}}(x) = \mathbf{w} \cdot \phi(x) \mid \mathbf{w} \in \mathbb{R}^d \}$$

All predictors

Hypothesis class

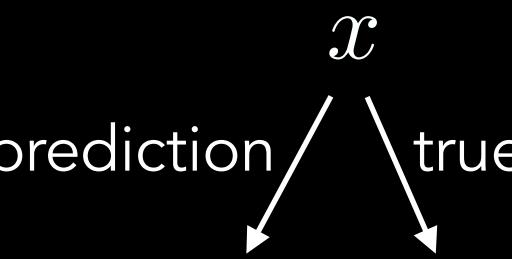
$$\star f_{\hat{\mathbf{w}}}$$

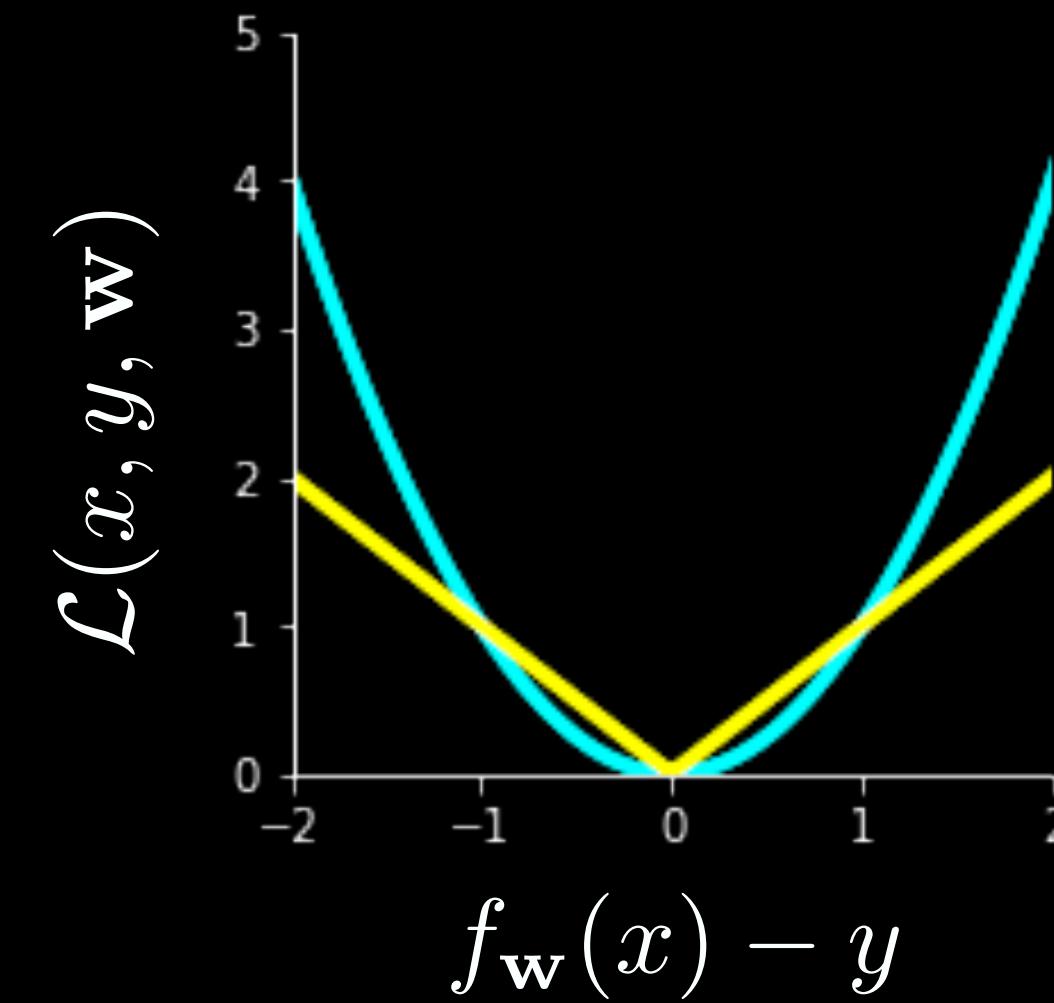
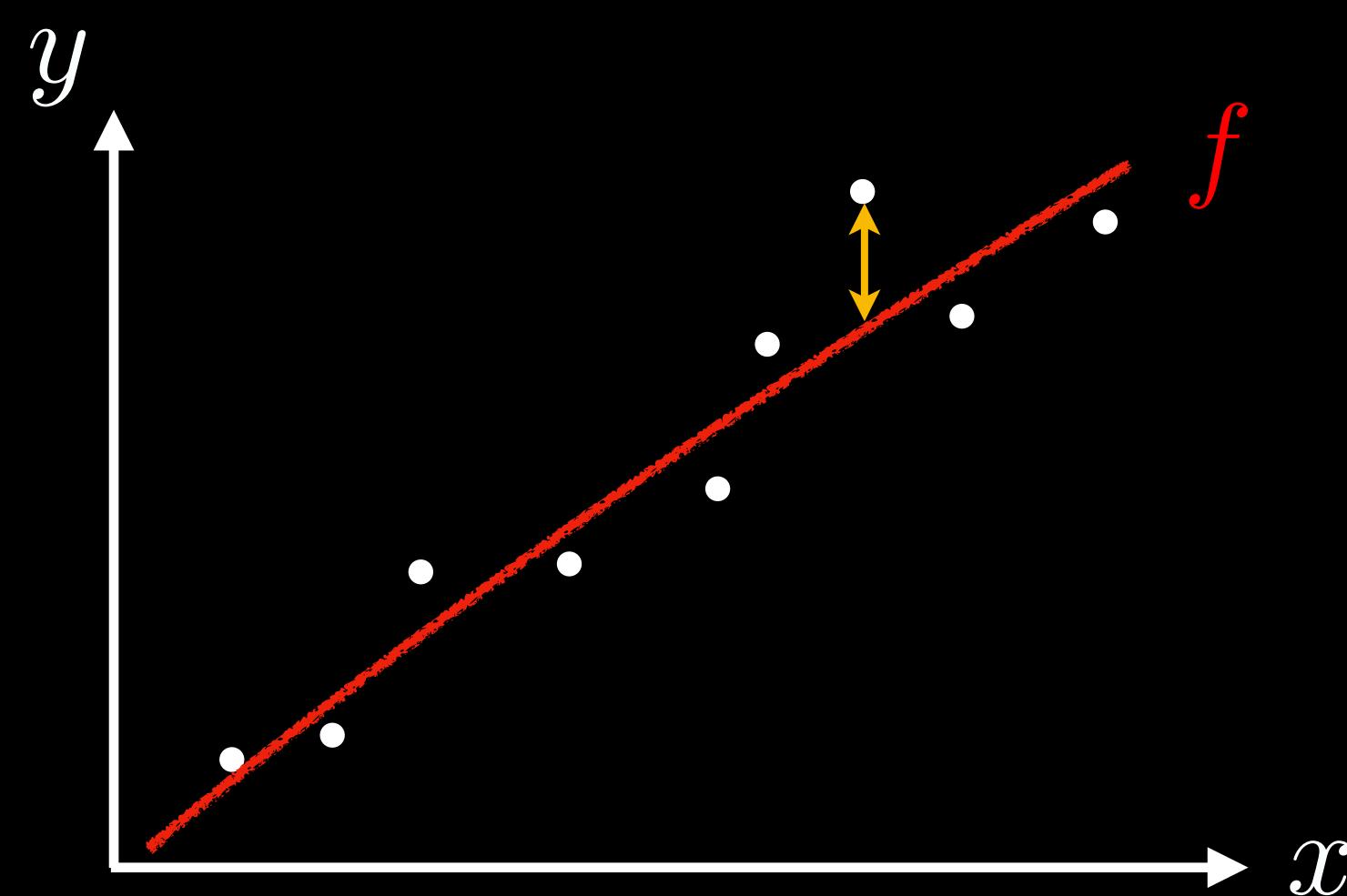
- Choose \mathcal{F} based on domain knowledge
- Find the best predictor $f_{\hat{\mathbf{w}}}(x)$ based on data

The loss function

A **loss function** quantifies how *bad* the predictor is

$$\mathcal{L}(x, y, \mathbf{w}) = \text{distance}(f_{\mathbf{w}}(x), y)$$





$$\mathcal{L}_{\text{sq}}(x, y, \mathbf{w}) = (f_{\mathbf{w}}(x) - y)^2$$

$$\mathcal{L}_{\text{abs}}(x, y, \mathbf{w}) = |f_{\mathbf{w}}(x) - y|$$

Minimize the loss

The **training loss** is the average loss over the data set

$$\mathcal{L}_{\text{train}}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \mathcal{L}(x, y, \mathbf{w})$$

objective

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\mathbf{w})$$

optimal predictor

$$f_{\hat{\mathbf{w}}}(x) = \hat{\mathbf{w}} \cdot \phi(x)$$

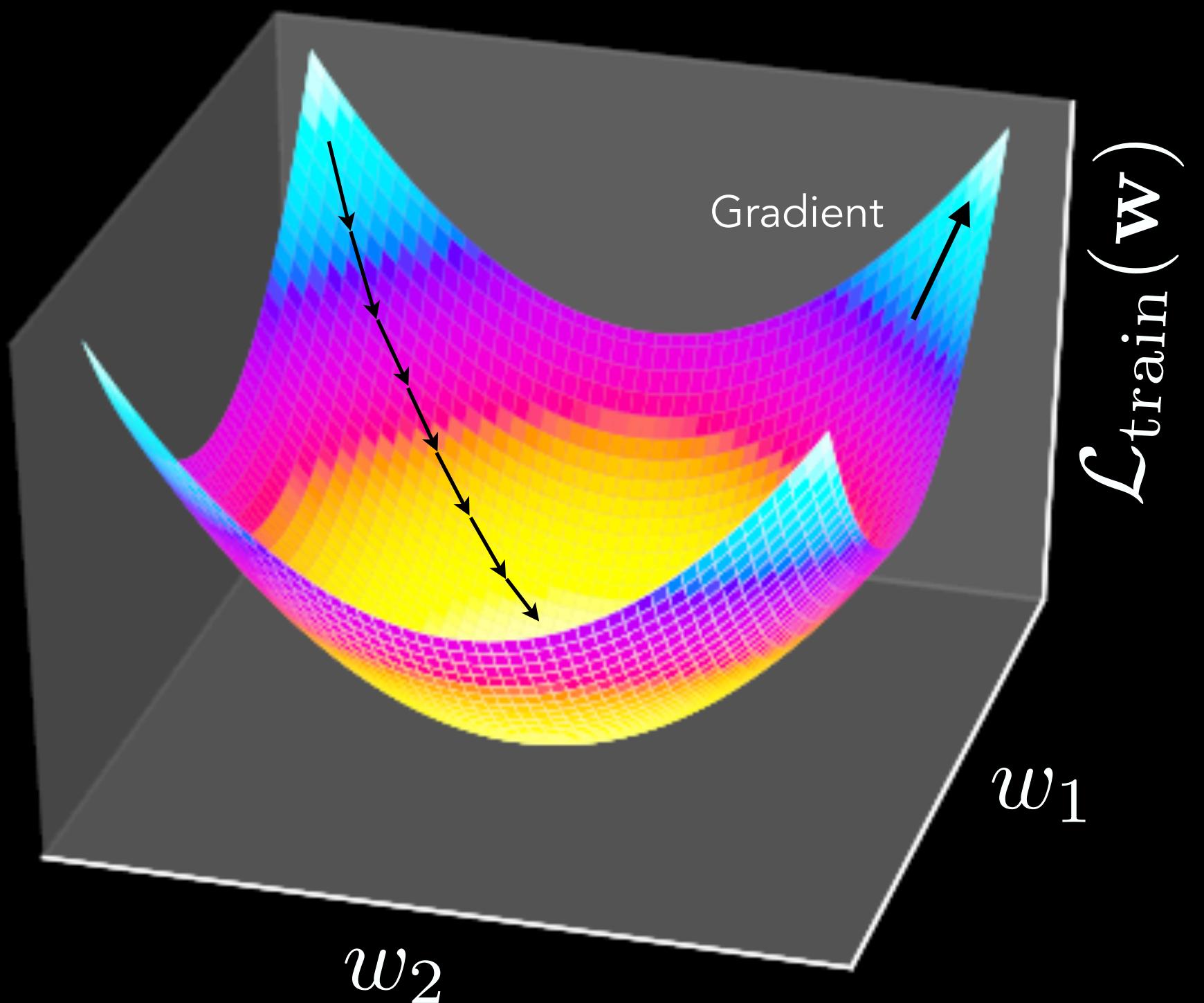
Gradient Descent

Training loss gradient

$$\nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \underbrace{2(\mathbf{w} \cdot \phi(x) - y)\phi(x)}_{\text{Prediction} - \text{True value}}$$

Gradient descent algorithm

```
epochs  
initialize  $\mathbf{w} = [0, \dots, 0]$ ;  
for  $t = 1, \dots, T$  do  
|  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}$   
end  
step size
```



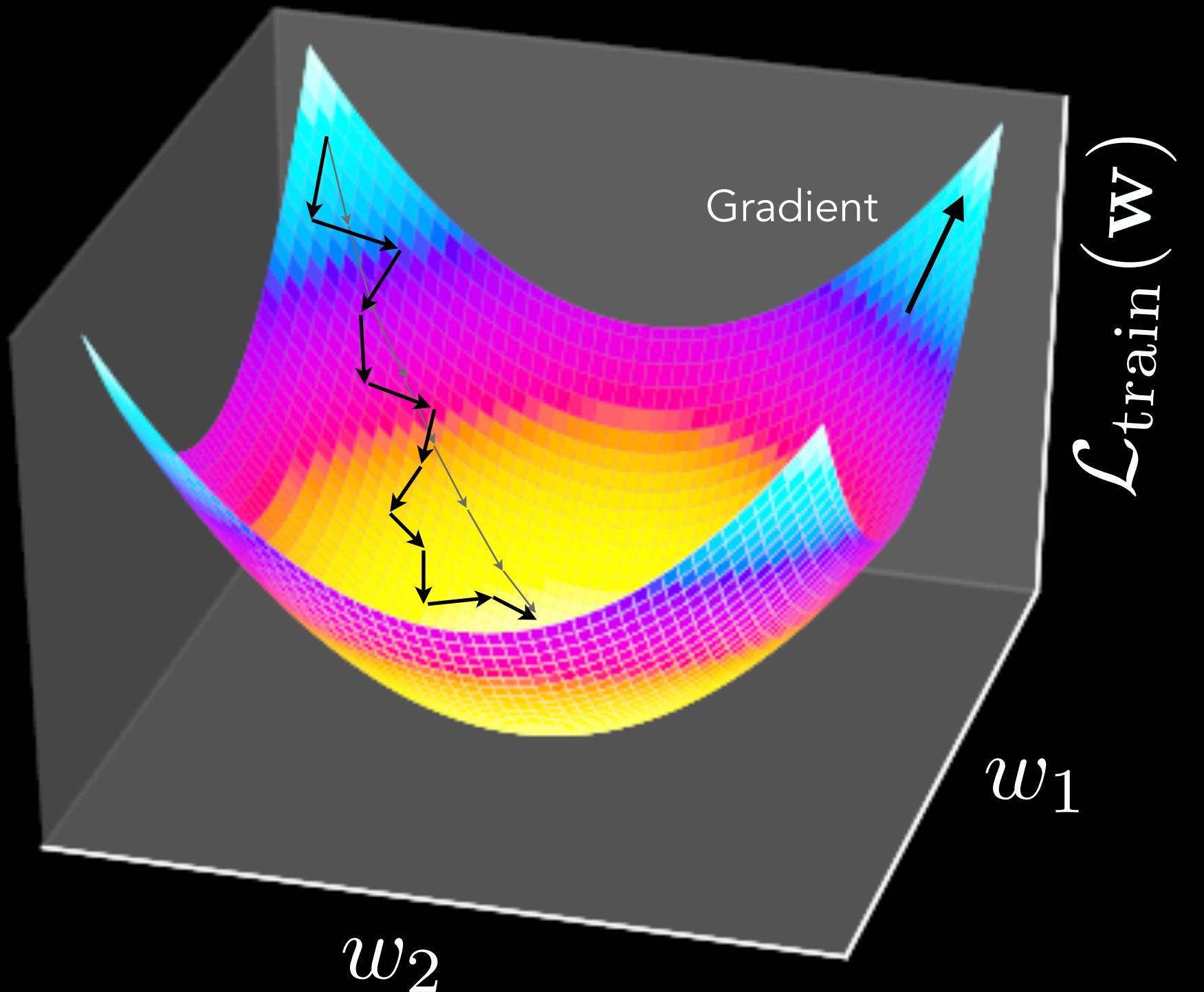
Stochastic Gradient Descent

Gradient descent is too expensive

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}$$

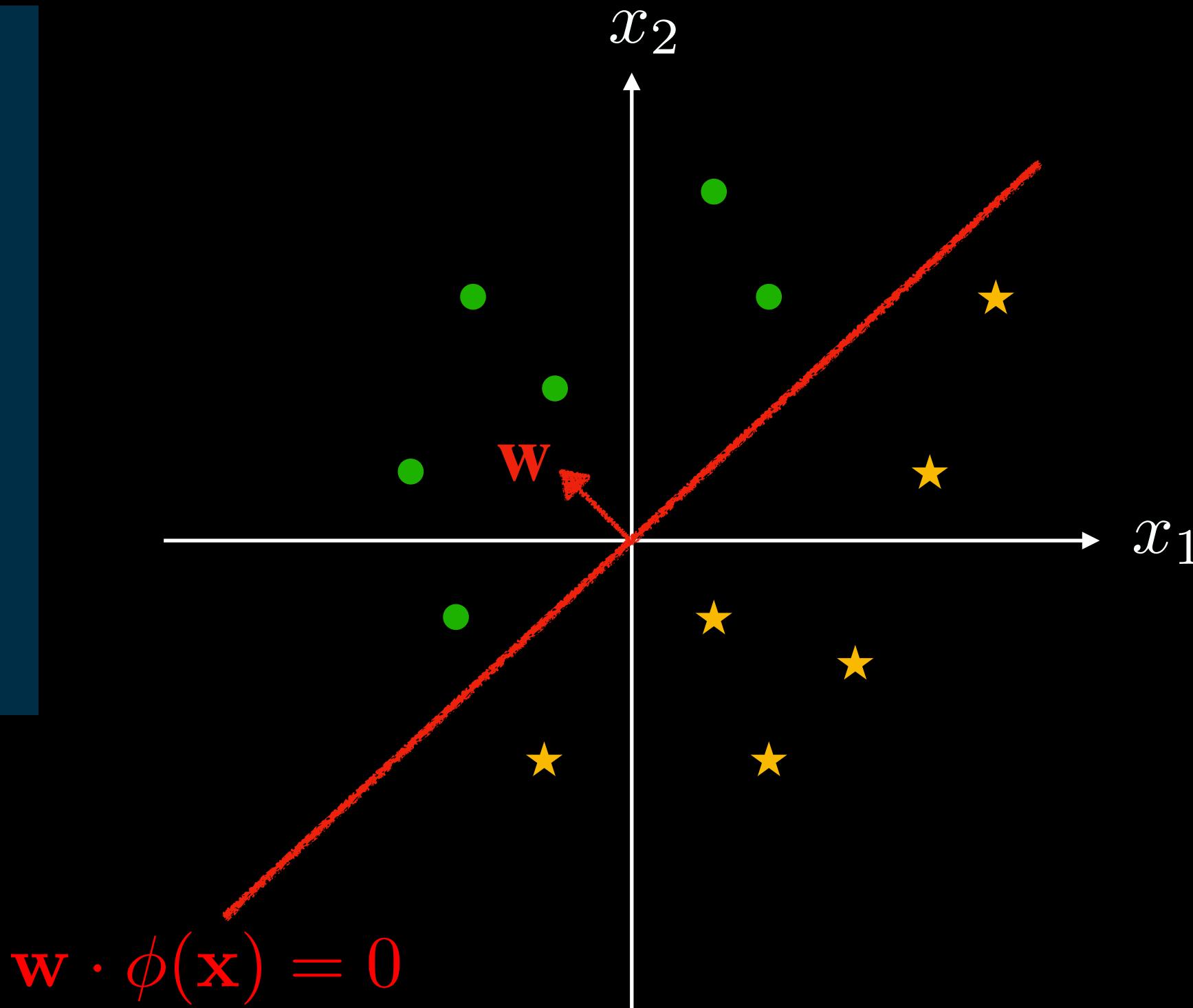
Stochastic gradient descent

```
initialize  $\mathbf{w} = [0, \dots, 0]$ ;
for  $t = 1, \dots, T$  do
    for  $(x, y) \in \mathcal{D}_{\text{train}}$  do
        |  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}(x, y, \mathbf{w})$ 
    end
end
```



Classification

x_1	x_2	y
-2	-1	•
3	1	★
2	3	•
1	-1	★
⋮		



$$y \in [-1, 1]$$

$$f_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \phi(\mathbf{x}))$$

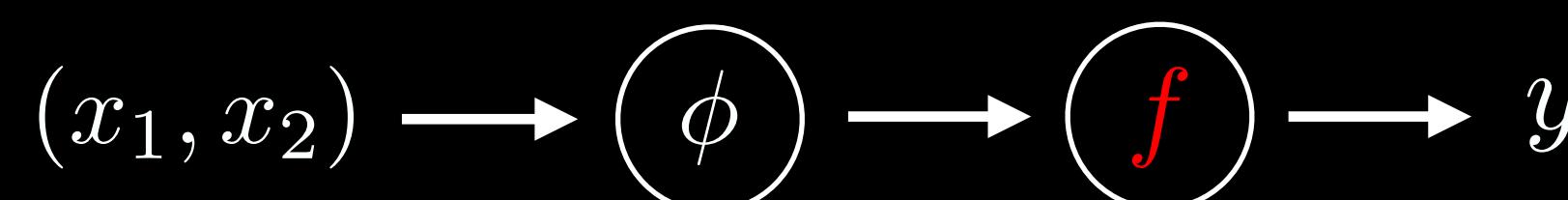
$$\mathcal{L}_{0-1} = \mathbf{1}[f_{\mathbf{w}}(\mathbf{x}) \neq y]$$

how confident?

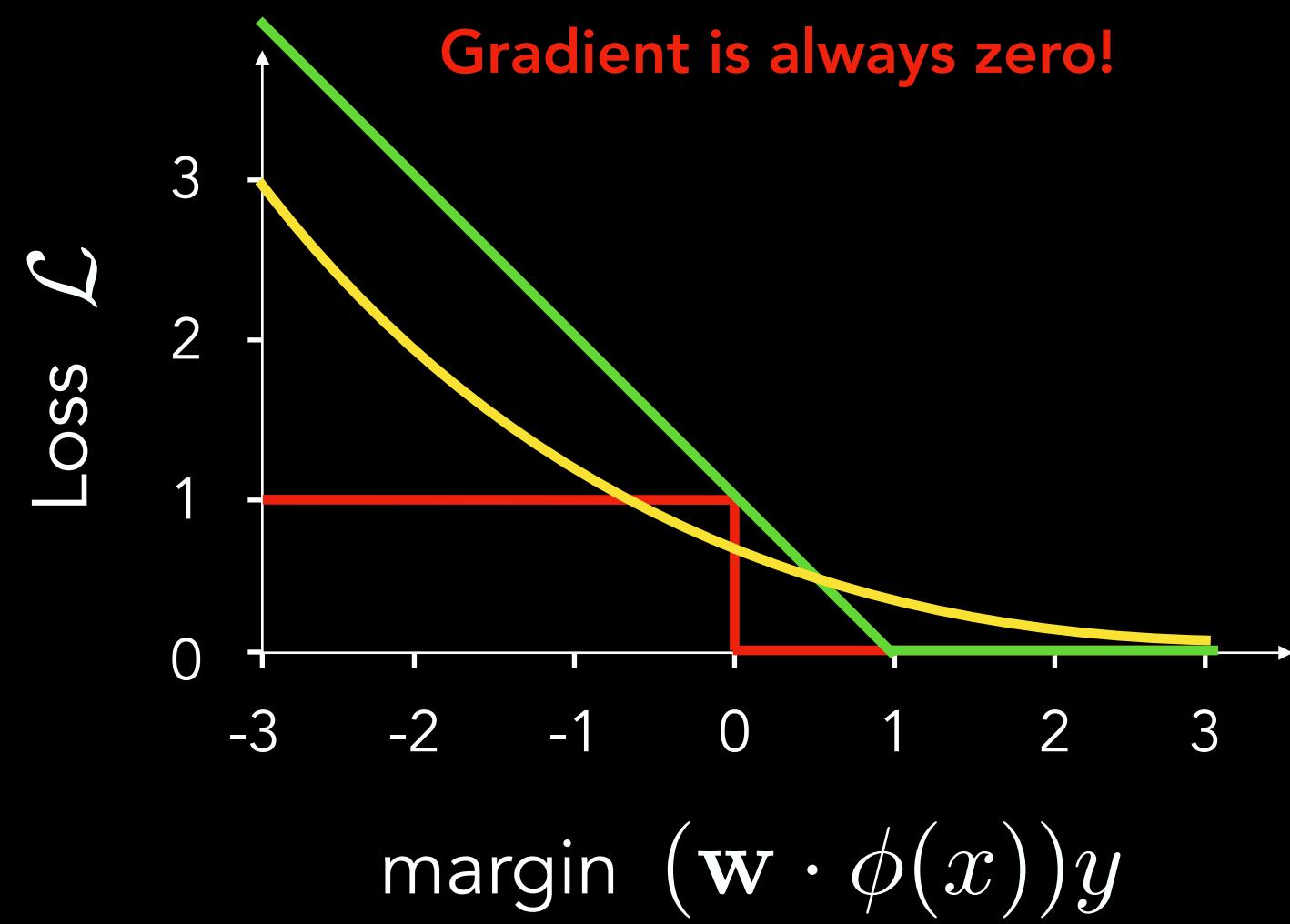
score
$\mathbf{w} \cdot \phi(\mathbf{x})$

how correct?

margin
$(\mathbf{w} \cdot \phi(\mathbf{x})) y$



Classification losses



$$\mathcal{L}_{0-1} = \mathbf{1} [(\mathbf{w} \cdot \phi(x)) y \leq 0]$$

$$\mathcal{L}_{\text{hinge}} = \max \{1 - (\mathbf{w} \cdot \phi(x)) y, 0\}$$

Support Vector Machines

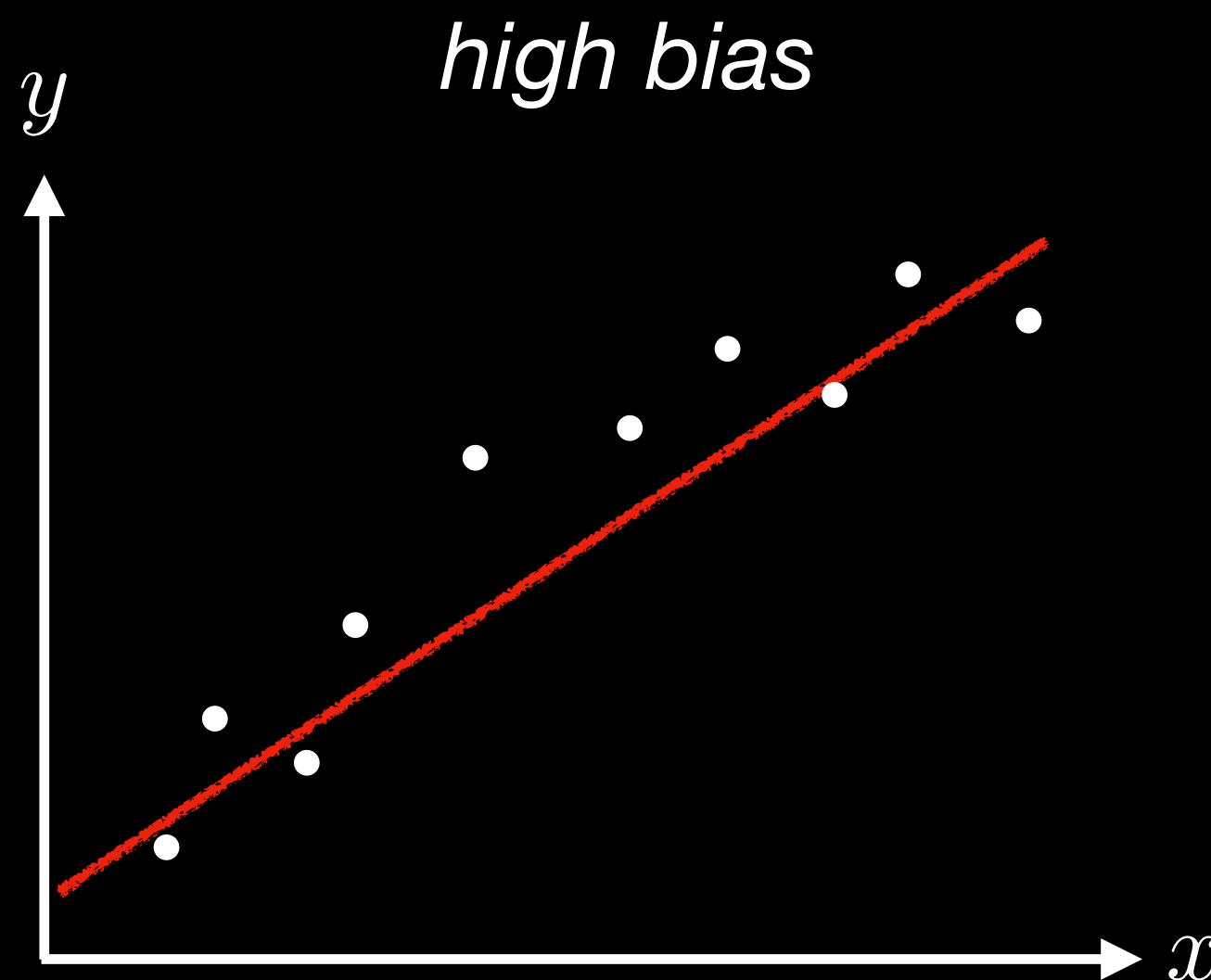
$$\mathcal{L}_{\text{logistic}} = \log \left(1 + e^{-(\mathbf{w} \cdot \phi(x)) y} \right)$$

Logistic Regression

Variance vs. Bias

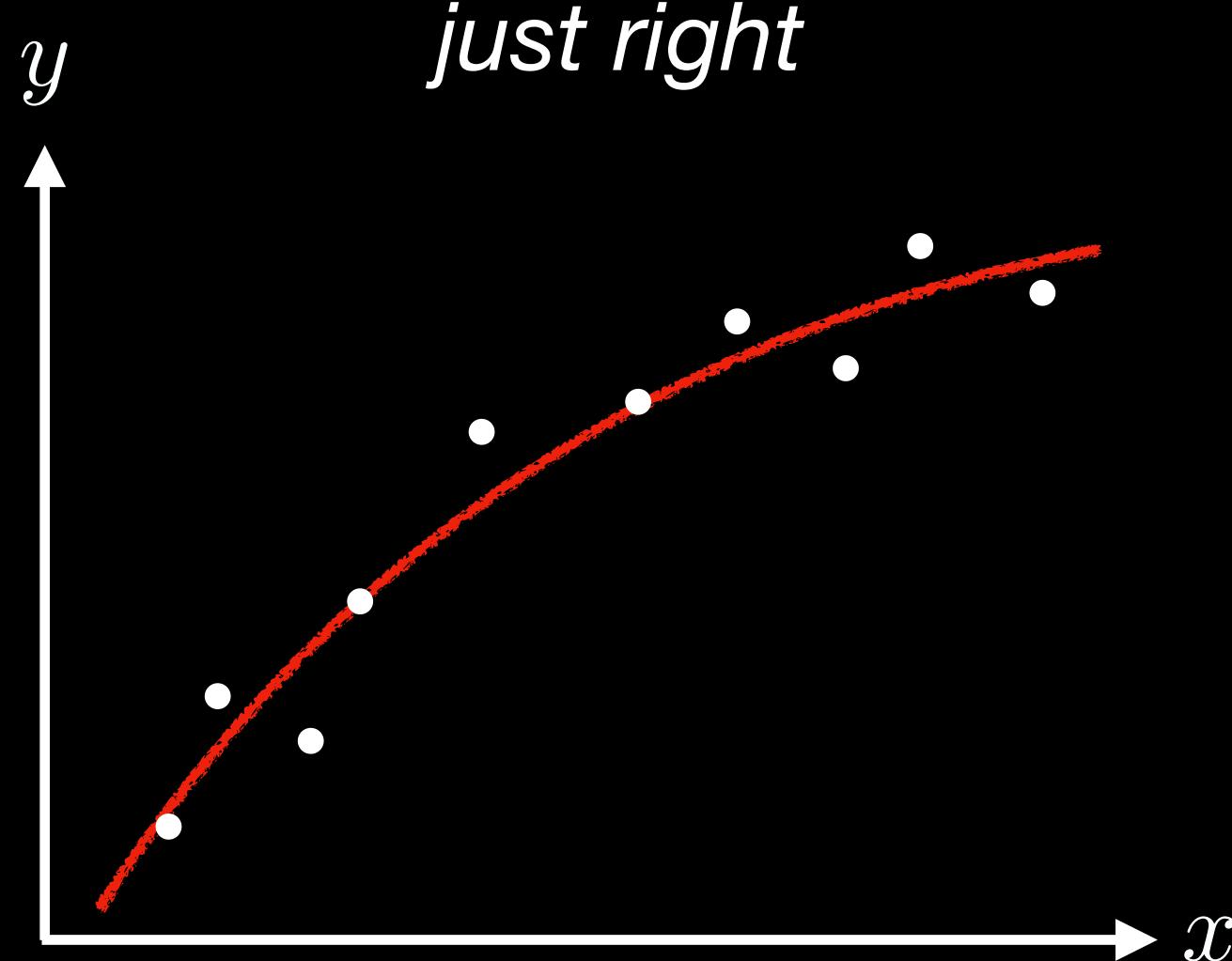
How do you choose $\phi(\cdot)$?

under-fitting



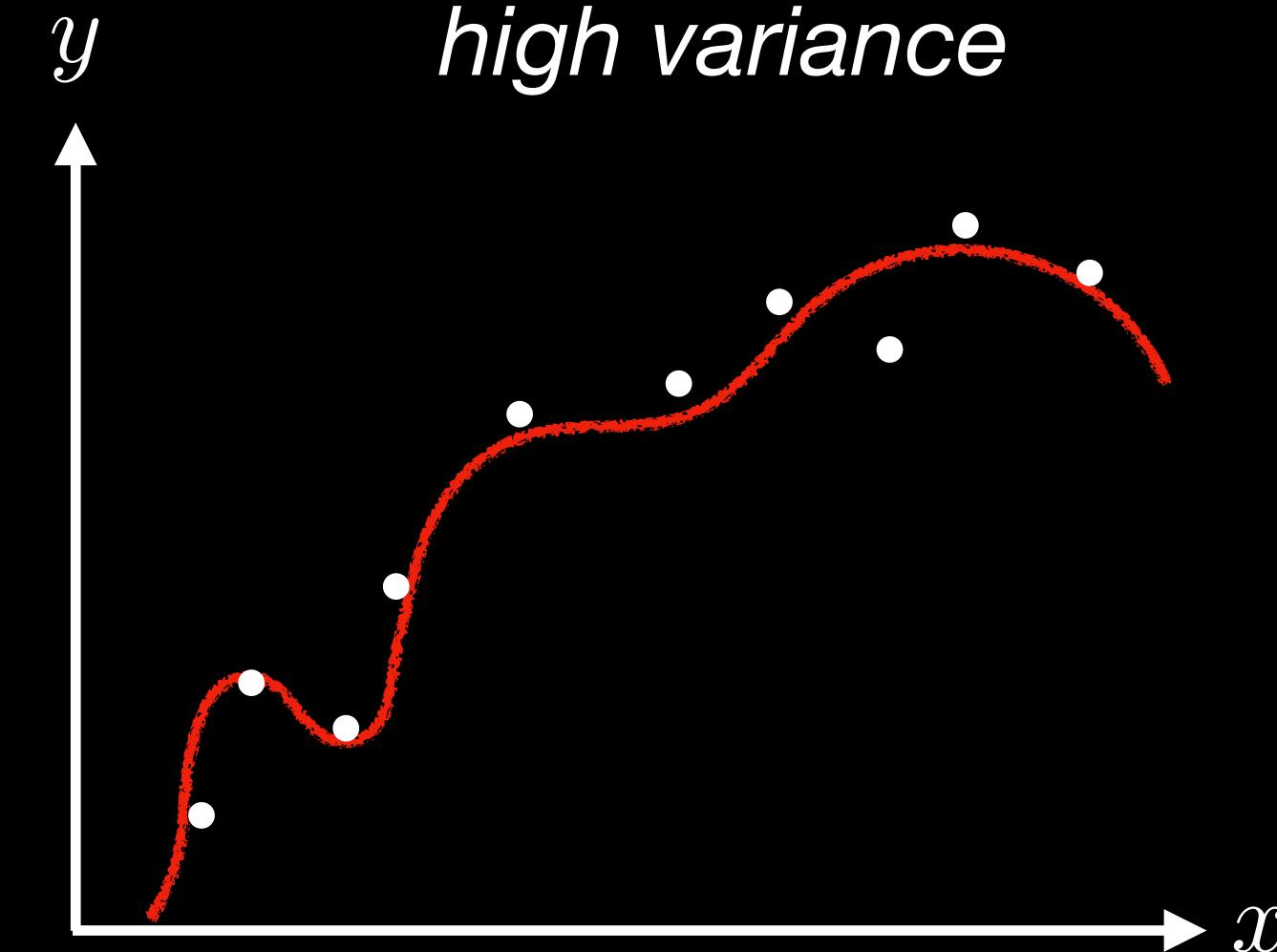
$$\phi(x) = [1, x]$$

just right



$$\phi(x) = [1, x, x^2]$$

over-fitting
high variance



$$\phi(x) = [1, x, x^2, \dots, x^n]$$

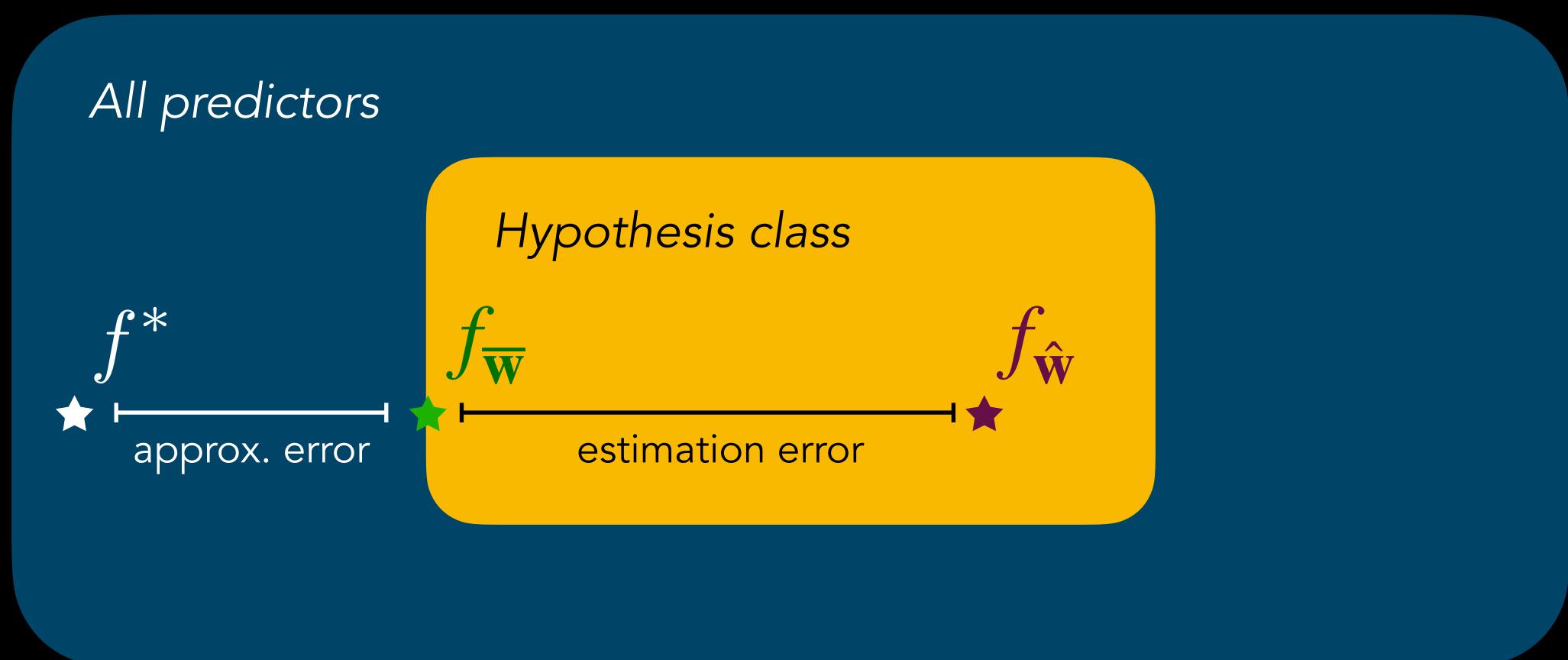
Two sources of error

Estimation error

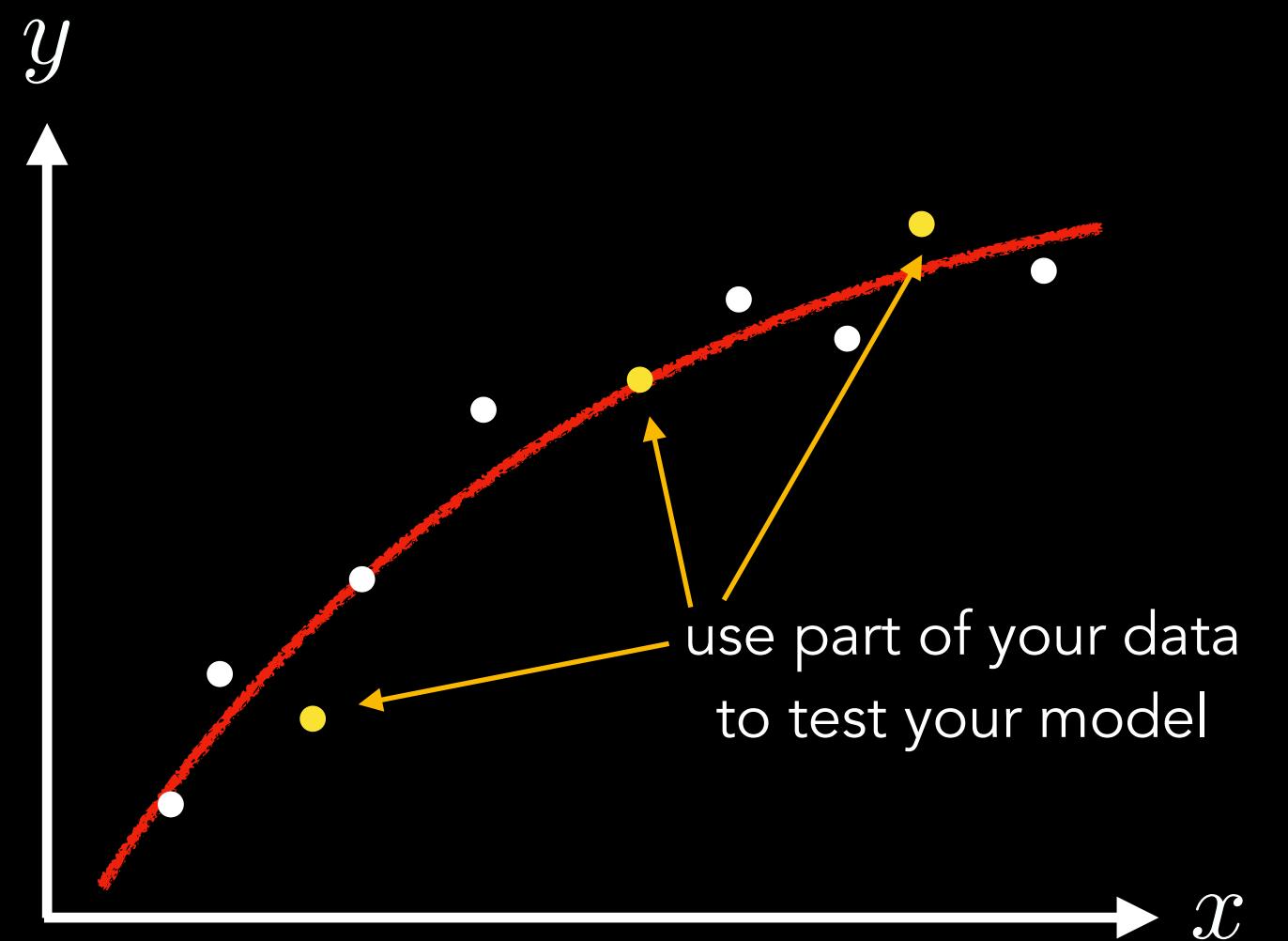
How good is the predictor, given the hypothesis class?

Approximation error

How good is the hypothesis class?

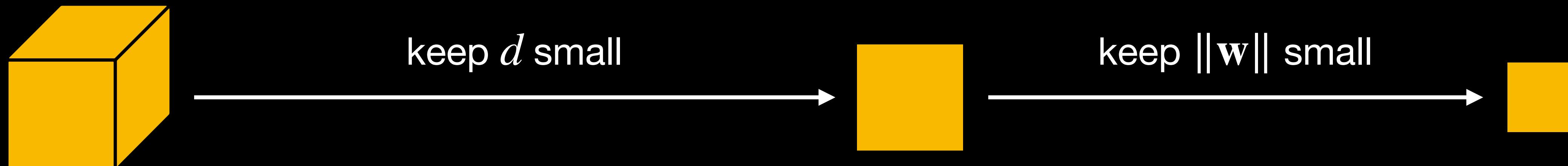


Use a test set



Designing the hypothesis class

$$\mathbf{w} \in \mathbb{R}^d$$



Remove features if they don't help

Use automatic feature selection

$$\mathcal{L} = \mathcal{L}_{\text{sq}} + \lambda \|\mathbf{w}\|_0$$

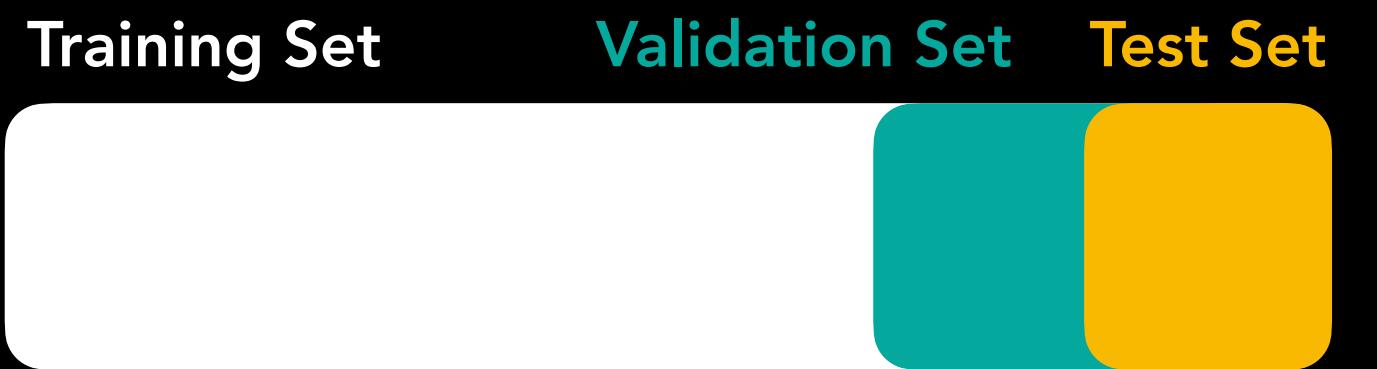
number of nonzero
elements

Use L_2 regularization
 $\mathcal{L} = \mathcal{L}_{\text{sq}} + \lambda \|\mathbf{w}\|^2$

Hyperparameters

How do you choose the

- Regularization parameter
- Number of iterations
- Steps size
- ...



The ML workflow

