

Applied Machine Learning

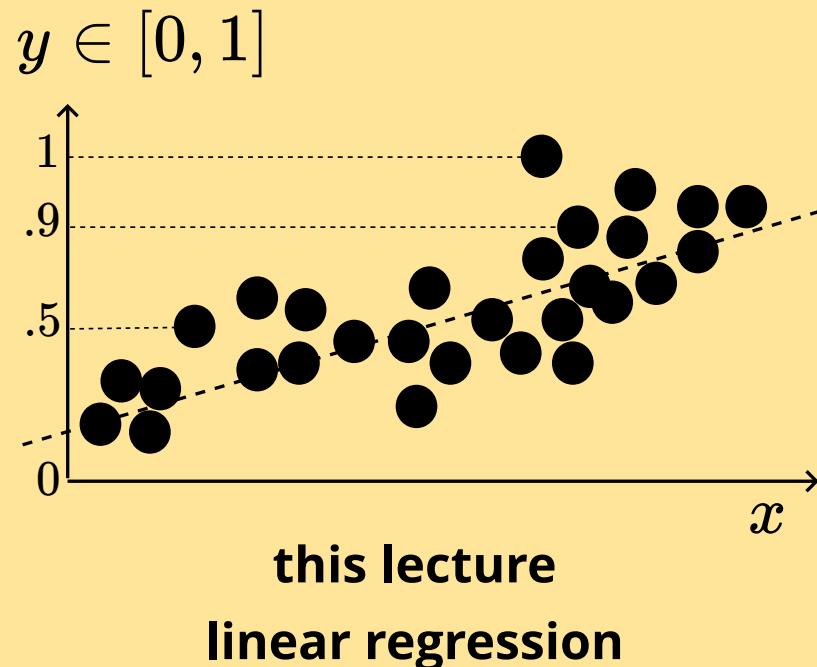
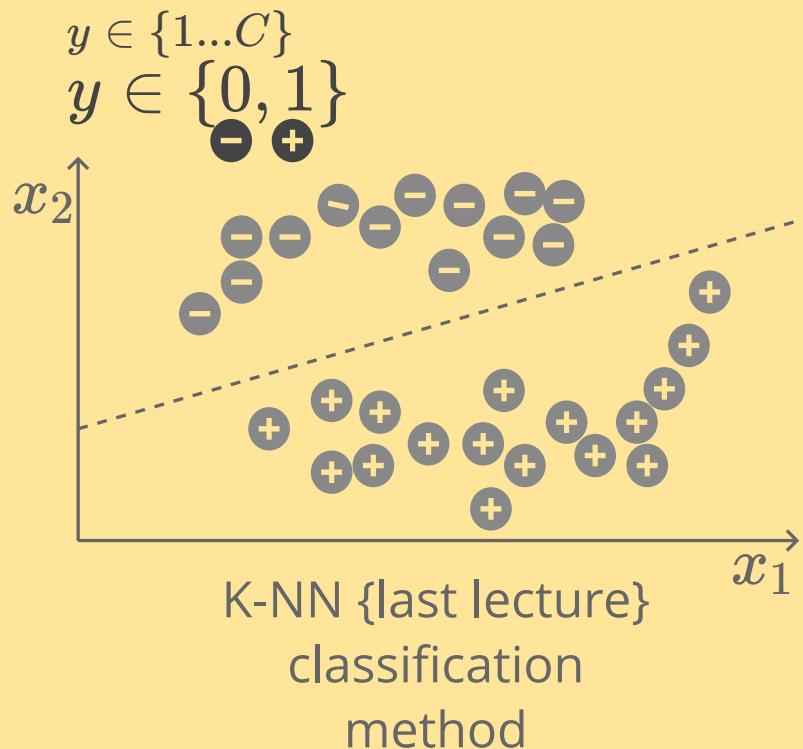
Linear Regression

Reihaneh Rabbany



COMP 551 (winter 2020)

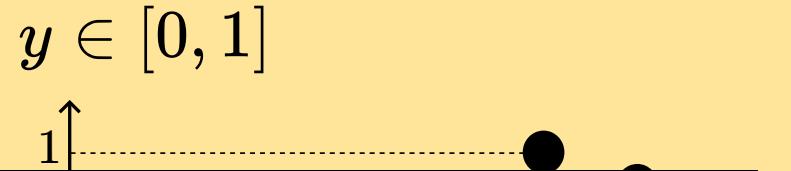
Linear Regression

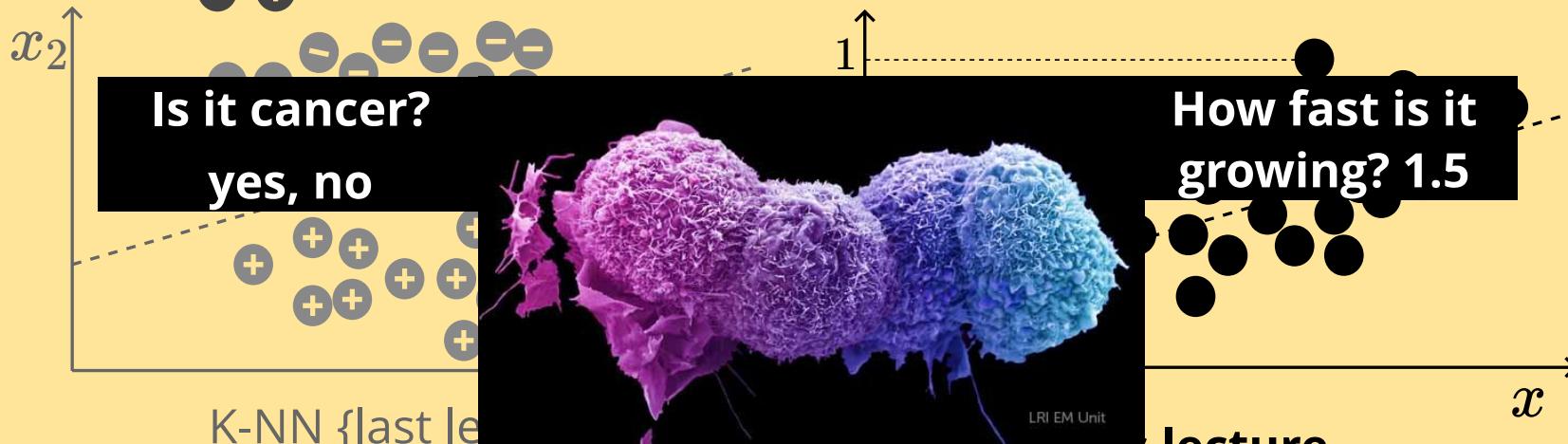


Regression

$$y \in \{1 \dots C\}$$

$$y \in \{0, 1\}$$


$$y \in [0, 1]$$




Regression: Example

Age-estimating
input: face
output: age

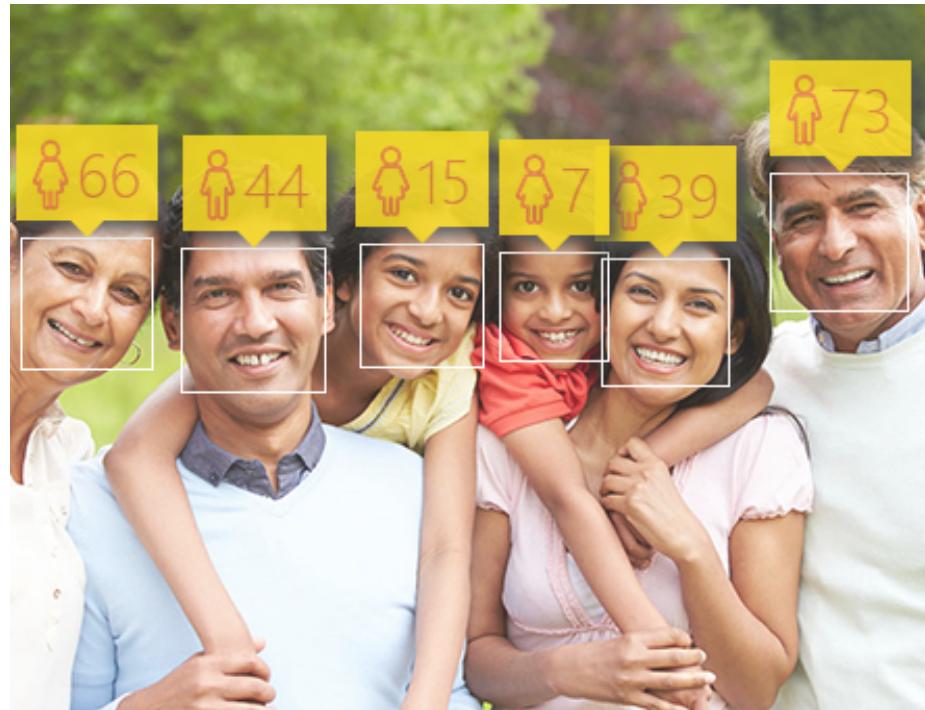


image from Microsoft age estimator here

Regression: Example

Protein folding
input: sequences
output: 3D structure

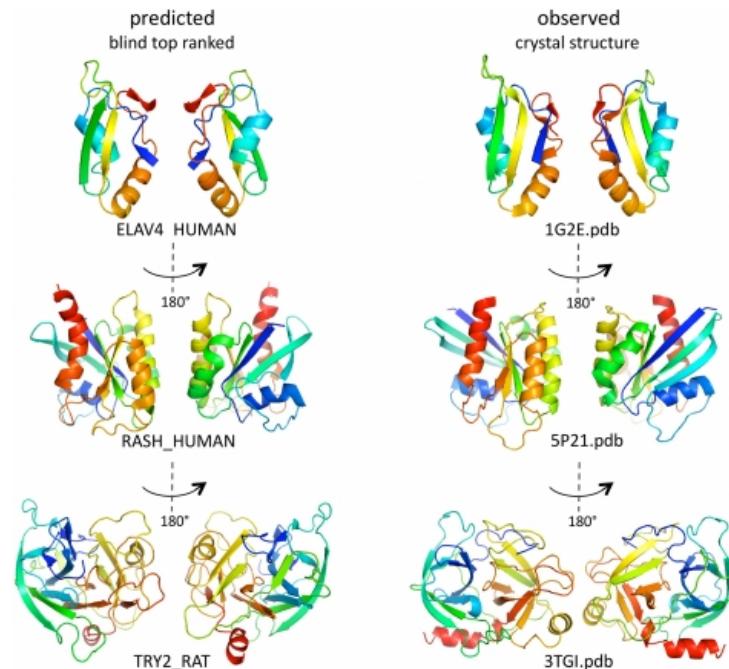


Image from Marks et al. [link](#)

Regression: Example

Colourization

input: gray scale image, output: colour image

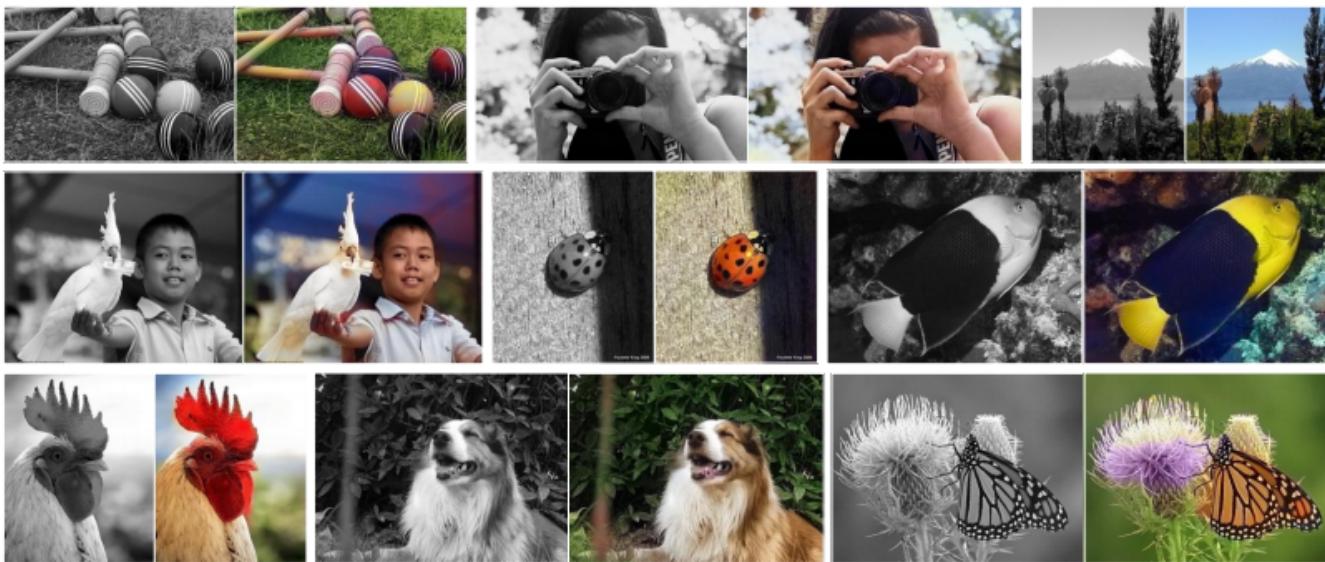
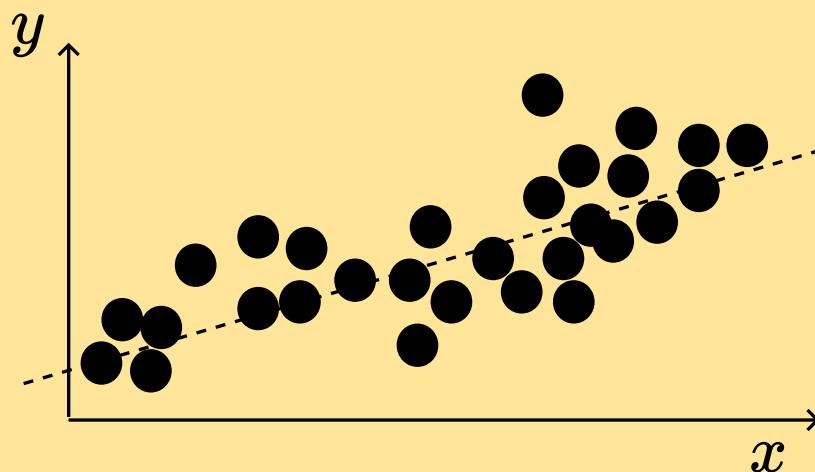


Image from Zhang et al. [link](#)

Learning objectives

linear least square regression

- Mathematical formulation
- Definition of loss function
- Optimization method
- Geometric interpretation
- ...



Motivation

History: method of least squares was invented by **Legendre** and **Gauss** (1800's)

Gauss used it to predict the future location of Ceres (largest asteroid in the asteroid belt)



ocean navigation

image from wiki history of navigation



Gauss
used it



Legendre
published it



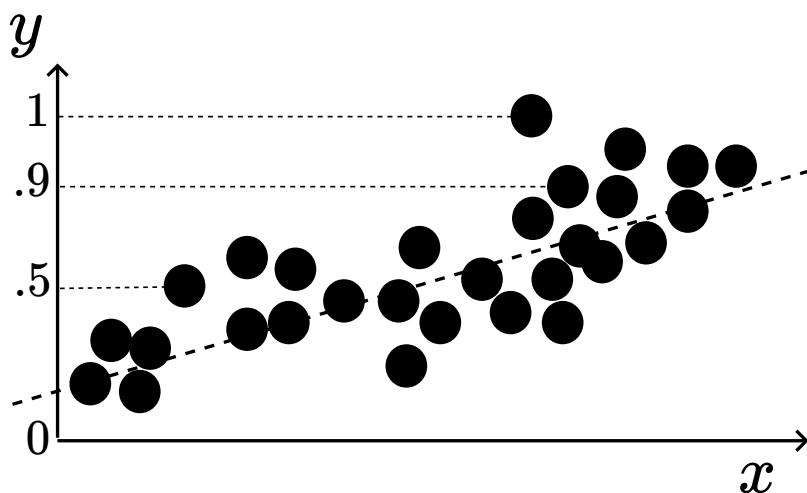
Pearson
named it regression

[read more on the history if interested](#)

Terminology

$x = [x_1, \dots, x_D]^T$ input, features, independent var., covariates

y target, dependent var., response var.

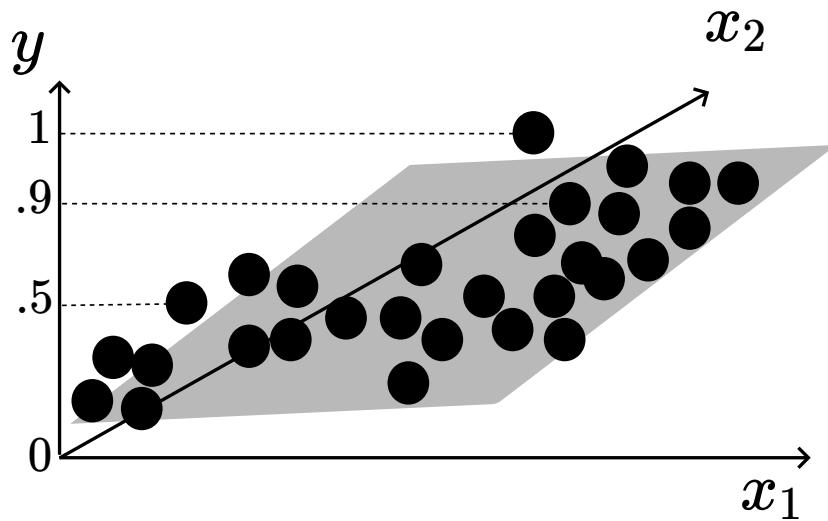


How many
dimensions are in
the input of this
example?

Terminology

$x = [x_1, \dots, x_D]^T$ input, features, independent var., covariates

y target, dependent var., response var.



Representing data

each instance: $x^{(n)} \in \mathbb{R}^D$ a (column) vector $x =$
 $y^{(n)} \in \mathbb{R}$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

a features

design matrix: concatenate all instances

- each row is a datapoint, each column is a feature

\mathbb{R} denotes set of real numbers



$$X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(N)T} \end{bmatrix} = \begin{bmatrix} x_1^{(1)}, & x_2^{(1)}, & \cdots, & x_D^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(N)}, & x_2^{(N)}, & \cdots, & x_D^{(N)} \end{bmatrix}$$

one instance

Representing data

design matrix: concatenate all instances

- each row is a datapoint, each column is a feature

instances

$$X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(N)T} \end{bmatrix} = \begin{bmatrix} x_1^{(1)}, & x_2^{(1)}, & \cdots, & x_D^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(N)}, & x_2^{(N)}, & \cdots, & x_D^{(N)} \end{bmatrix} \begin{matrix} \text{features} \\ \in \mathbb{R}^{N \times D} \end{matrix} \quad Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

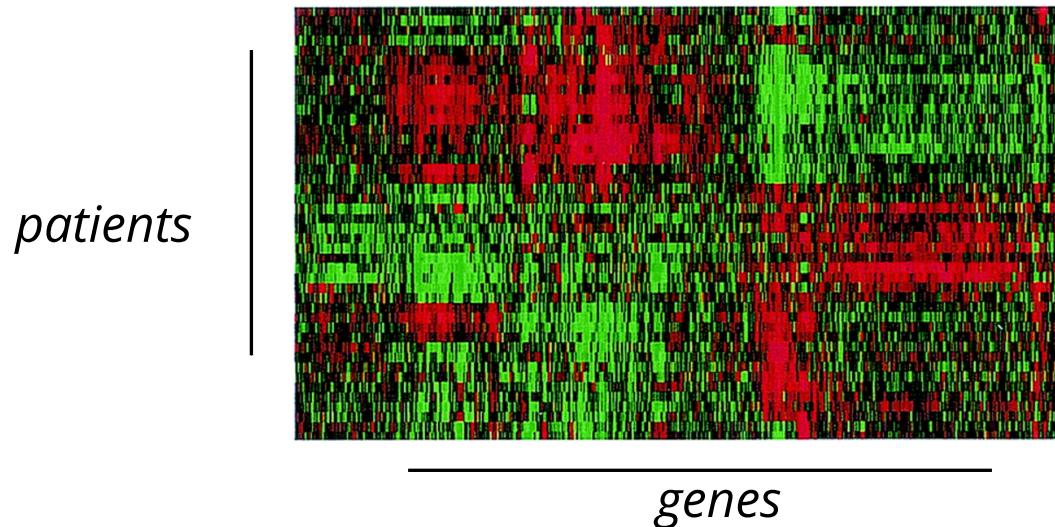
$$\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$$

we assume **N** instances in the dataset
each instance has **D** features indexed by **d**

Example: representing data

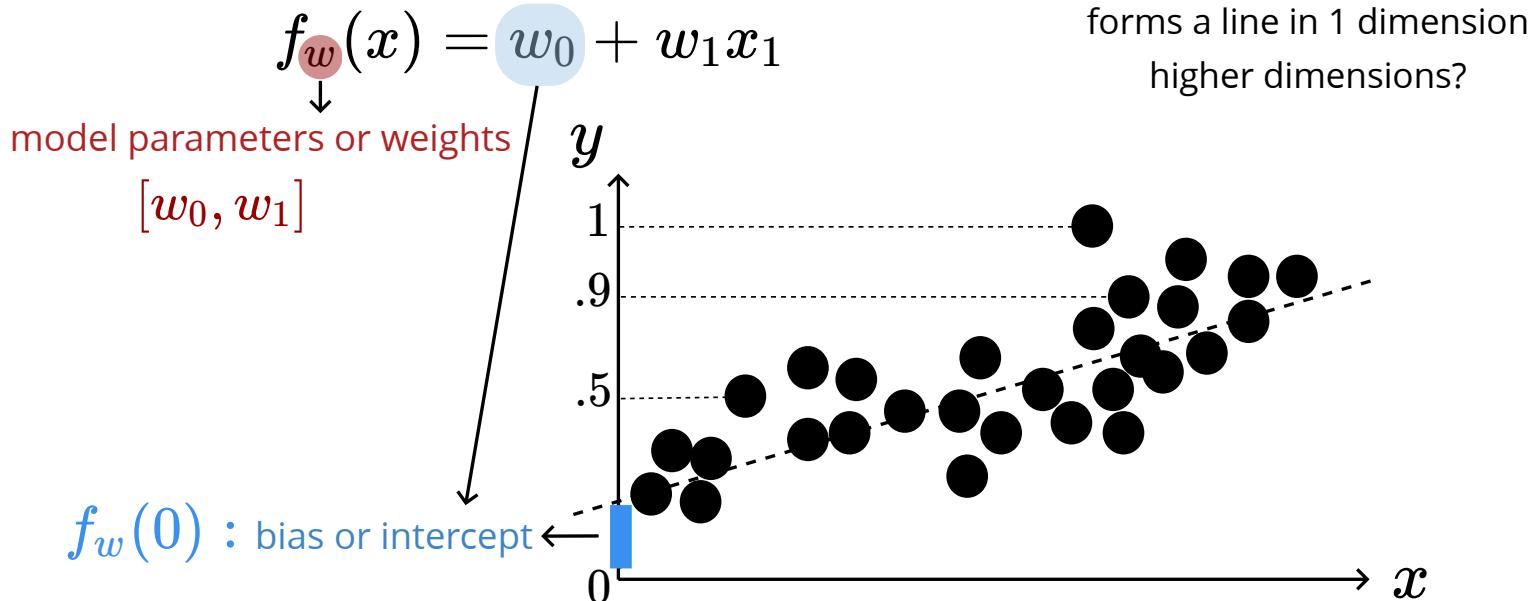
design matrix: each row is a datapoint, each column is a feature

Micro array data (X), contains gene expression level
labels (y) can be {cancer/no cancer} label for each patient



Linear model $D = 1$

assuming a scalar output $f_w : \mathbb{R} \rightarrow \mathbb{R}$
can generalize to a vector



Linear model

assuming a scalar output $f_w : \mathbb{R}^D \rightarrow \mathbb{R}$
can generalize to a vector

$$f_w(x) = w_0 + w_1 x_1 + \dots + w_D x_D$$

model parameters or weights
bias or intercept



```
yh_n = np.dot(w, x_n)
```

simplification

$$\text{concatenate a 1 to } x \rightarrow x = [1, x_1, \dots, x_D]^T \quad \triangleright \quad f_w(x) = w^T x$$
$$w = [w_0, w_1, \dots, w_D]^T$$

How should we pick the weights?

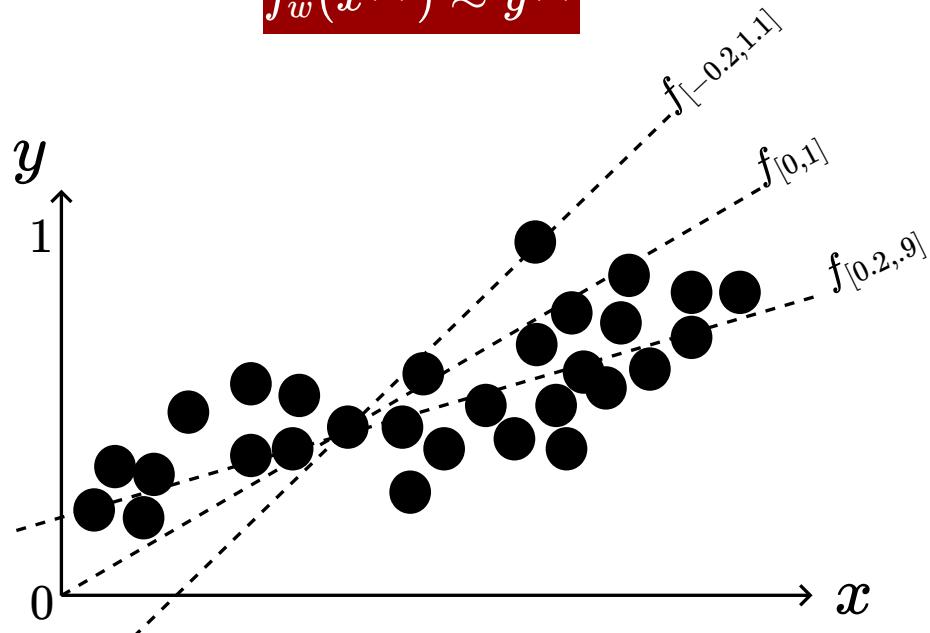
Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)} \quad \forall n$

$$f_w(x^{(n)}) \approx y^{(n)}$$

$$f_w(x) = w^T x$$

$$w = [w_0, w_1]$$

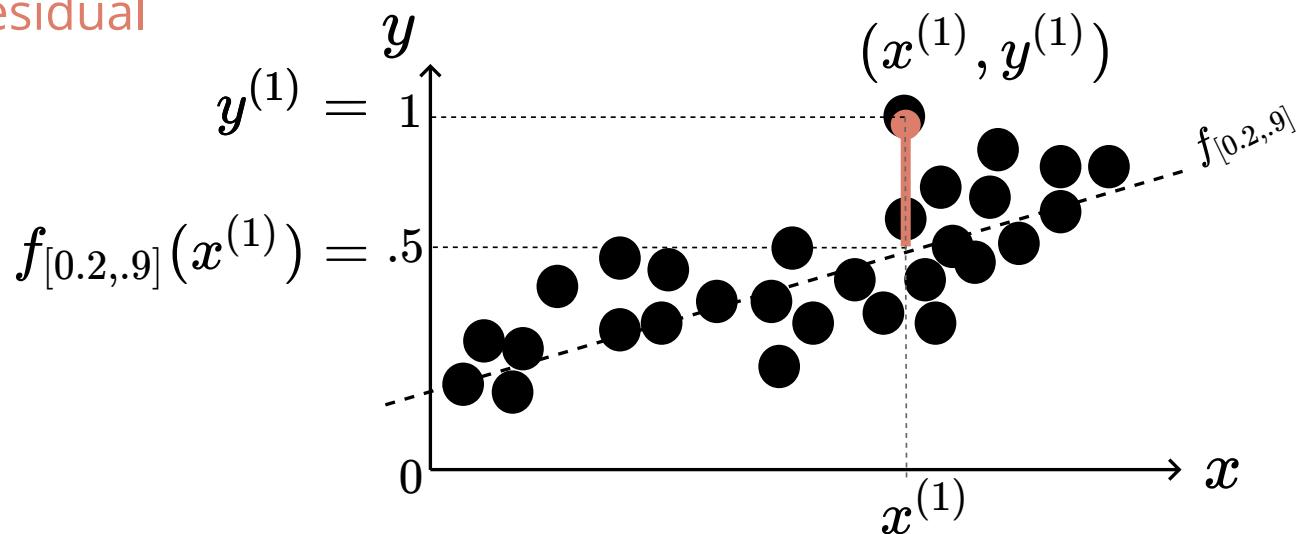


Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)} \quad \forall n$

compute
residual

$$f_w(x^{(n)}) \approx y^{(n)}$$



Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)}$ $\forall n$

$$\hat{y}^{(n)} = f_w(x^{(n)}) \approx y^{(n)}$$

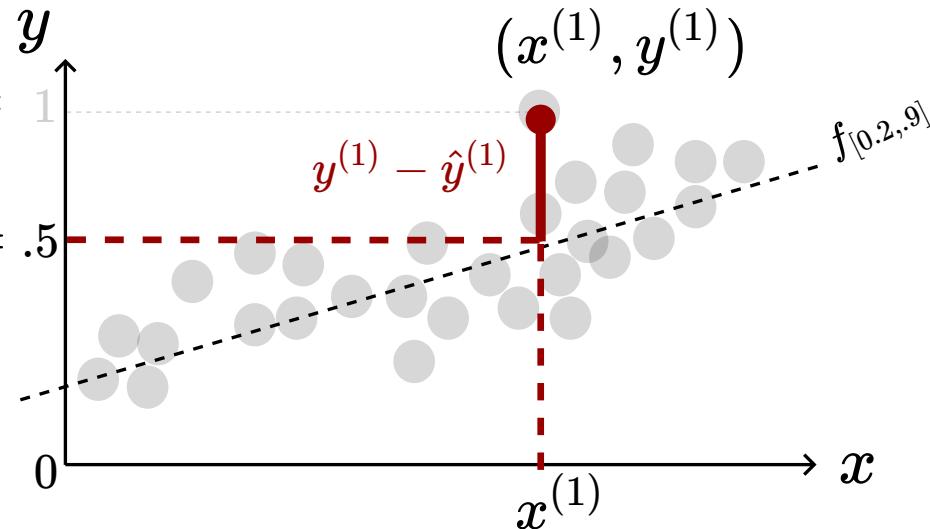
why residual?

true:

$$y^{(1)} = 1$$

predicted:

$$\hat{y}^{(1)} = f(x^{(1)}) = .5$$

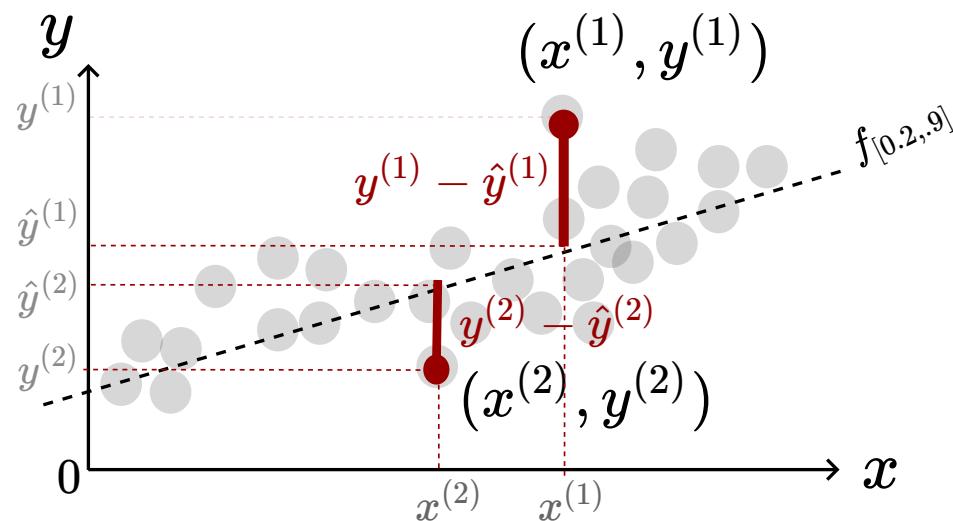


Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)}$ $\forall n$

$$f_w(x^{(n)}) \approx y^{(n)}$$

how to sum all residuals?



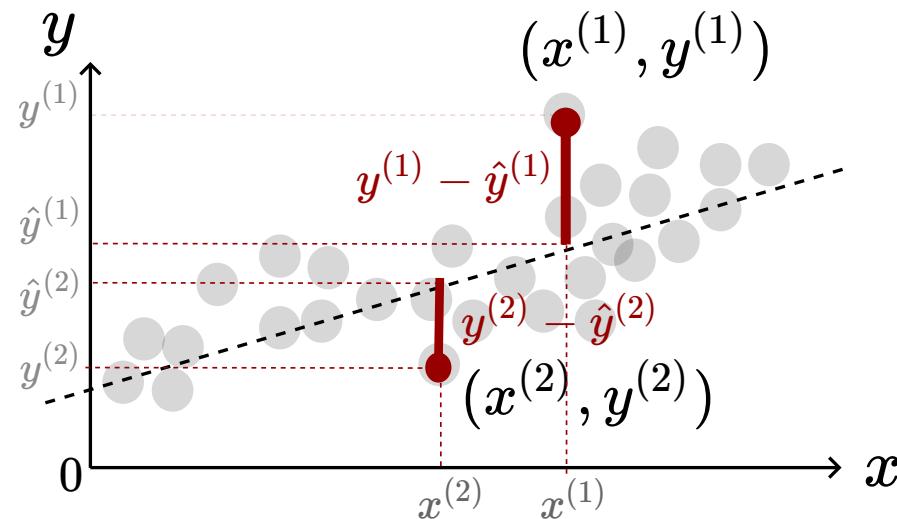
Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)} \quad \forall n$

$$f_w(x^{(n)}) \approx y^{(n)}$$

square error **loss**
(a.k.a. **L2** loss)

$$L(y, \hat{y}) \triangleq (y - \hat{y})^2$$



Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)}$ $\forall n$

$$f_w(x^{(n)}) \approx y^{(n)}$$

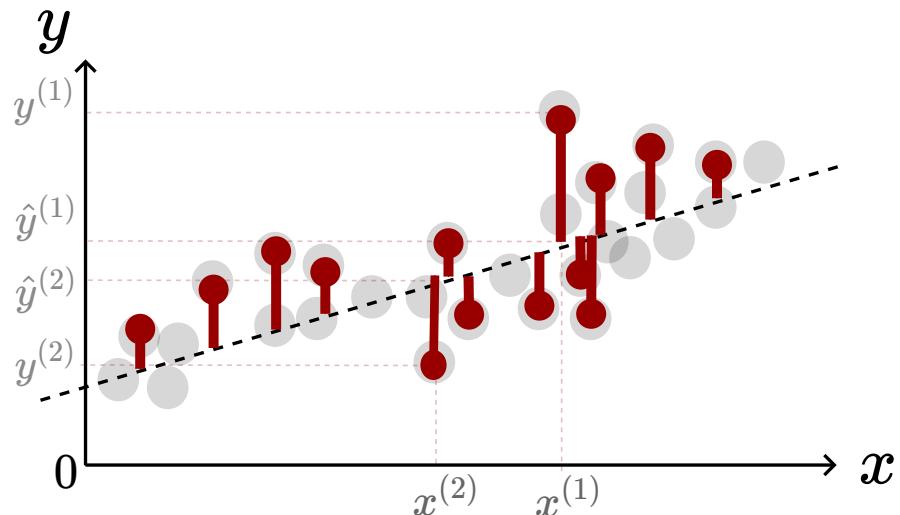
linear least squares

cost function

$$J(w) = \frac{1}{2} \sum_n \left(y^{(n)} - w^T x^{(n)} \right)^2$$

mean squared error

$$\frac{1}{N} J(w)$$



Loss function

objective: find parameters to **fit the data** $x^{(n)}, y^{(n)} \quad \forall n$

$$f_w(x^{(n)}) \approx y^{(n)}$$

minimize a measure of difference between $\hat{y}^{(n)} = f_w(x^{(n)})$ and $y^{(n)}$

square error **loss** (a.k.a. **L2 loss**) $L(y, \hat{y}) \triangleq \frac{1}{2}(y - \hat{y})^2$

for a single instance

versus

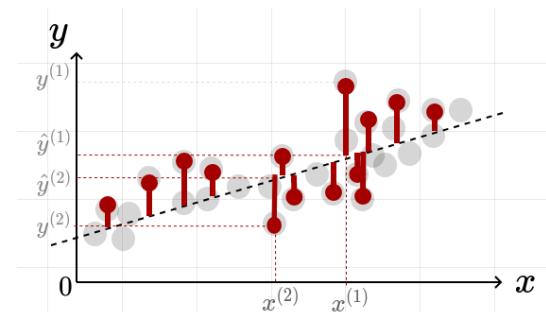
for future convenience

for the whole dataset

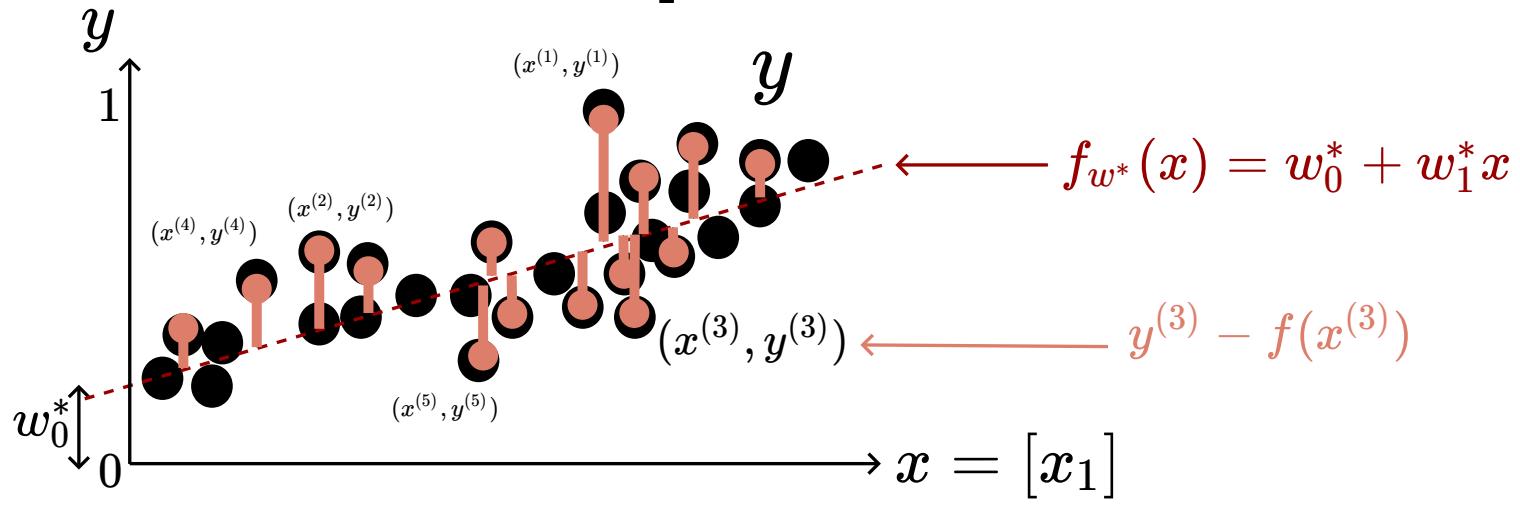
linear least squares **cost function**

$$J(w) = \frac{1}{2} \sum_n \left(y^{(n)} - w^T x^{(n)} \right)^2$$

$$w^* = \arg \min_w J(w)$$



Example (D=1)



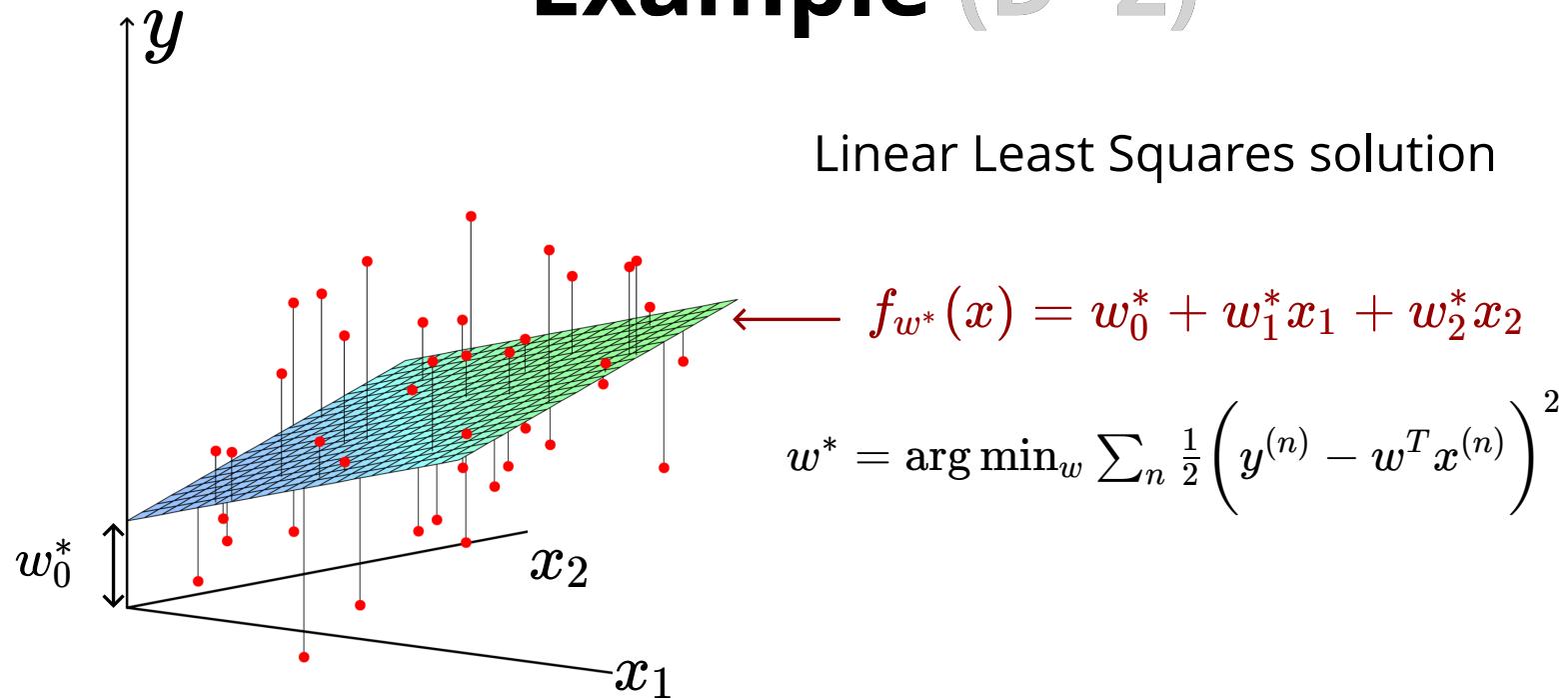
Linear Least Squares

cost:

$$J(w) = \sum_n \frac{1}{2} \left(y^{(n)} - w^T x^{(n)} \right)^2$$

solution: $w^* = \arg \min_w \sum_n \frac{1}{2} \left(y^{(n)} - w^T x^{(n)} \right)^2$

Example (D=2)



Minimizing the cost

Simple case: $D = 1$

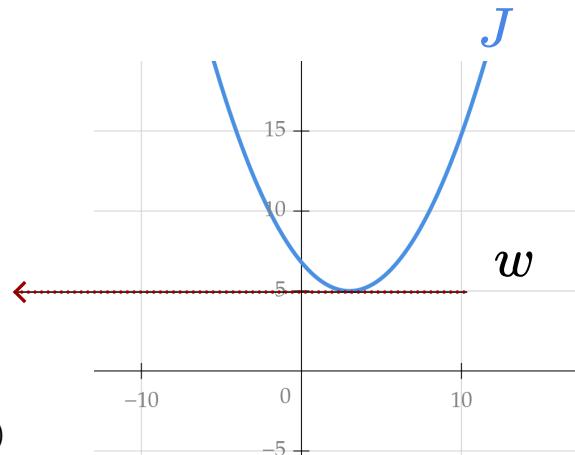
model $f_w(x) = wx$
both scalar

cost function $J(w) = \frac{1}{2} \sum_n (y^{(n)} - wx^{(n)})^2$

derivative $\frac{dJ}{dw} = \sum_n x^{(n)} (wx^{(n)} - y^{(n)})$

setting the derivative to zero $\rightarrow w^* = \frac{\sum_n x^{(n)} y^{(n)}}{\sum_n x^{(n)} x^{(n)}}$

global minimum because cost is smooth and convex



Minimizing the cost

for a multivariate function $J(w_0, w_1)$

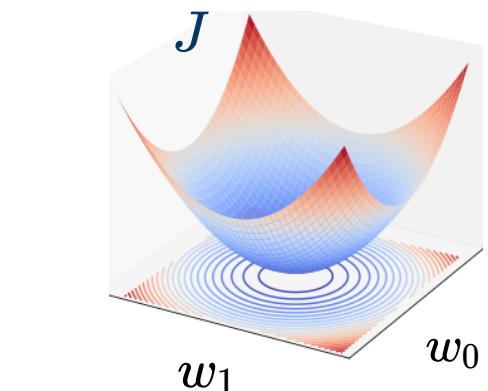
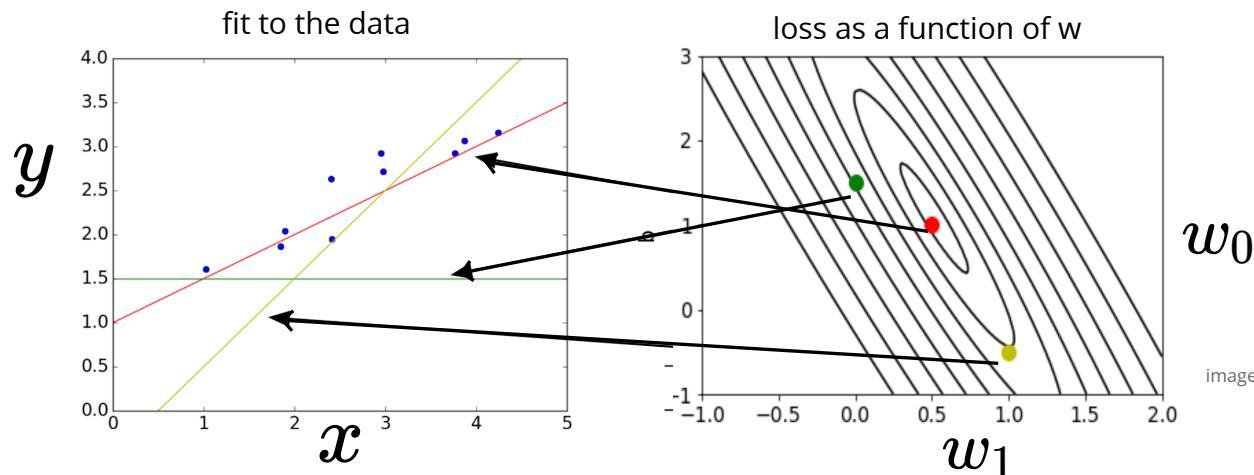


image: Grosse, Farahmand, Carrasquilla

The objective is a smooth function of w

Find minimum by setting partial derivatives to zero

image: Grosse, Farahmand, Carrasquilla

Gradient

for a multivariate function $J(w_0, w_1)$

partial derivatives instead of derivative

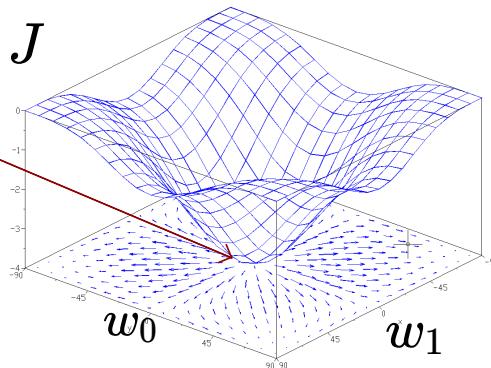
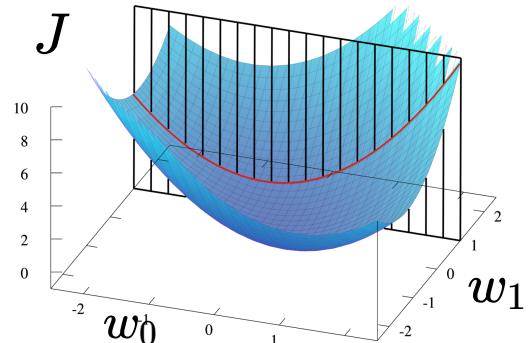
= derivative when other vars. are fixed

$$\frac{\partial}{\partial w_1} J(w_0, w_1) \triangleq \lim_{\epsilon \rightarrow 0} \frac{J(w_0, w_1 + \epsilon) - J(w_0, w_1)}{\epsilon}$$

critical point: all partial derivatives are zero

gradient: vector of all partial derivatives

$$\nabla J(w) = [\frac{\partial}{\partial w_1} J(w), \dots, \frac{\partial}{\partial w_D} J(w)]^T$$



Finding w (any D)

$$f_w(x) = w^T x$$

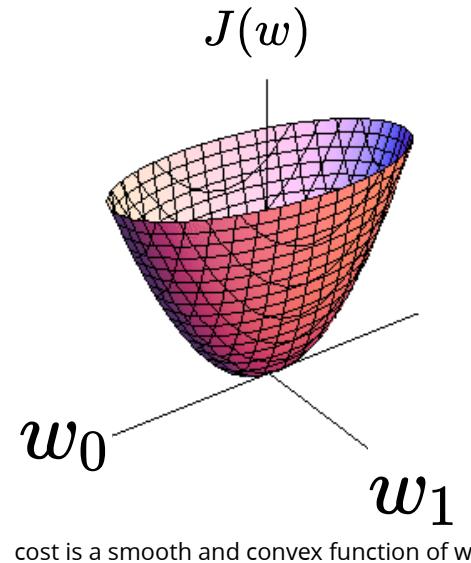
setting $\frac{\partial}{\partial w_i} J = 0$

$$\frac{\partial}{\partial w_i} \sum_n \frac{1}{2} (y^{(n)} - f_w(x^{(n)}))^2 = 0$$

using **chain rule**: $\frac{\partial J}{\partial w_i} = \frac{dJ}{df_w} \frac{\partial f_w}{\partial w_i}$

we get $\sum_n (w^T x^{(n)} - y^{(n)}) x_d^{(n)} = 0 \quad \forall d \in \{1, \dots, D\}$

D equations with D unknowns



cost is a smooth and convex function of w

Matrix form

instead of $\hat{y}^{(n)} = \mathbf{w}^T \mathbf{x}^{(n)}$
 $\in \mathbb{R}$ $1 \times D$ $D \times 1$

D is in fact dimensions of the input +1
due to the simplification and adding the
bias/intercept term

use **design matrix** to write $\hat{\mathbf{y}} = \mathbf{X} \mathbf{w}$
 $N \times 1$ $N \times D$ $D \times 1$

$$\hat{y}^{(1)} = w_0 + x_1^{(1)}w_1 + x_2^{(1)}w_2 + \dots + x_D^{(1)}w_D$$

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(N)} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)}, & x_2^{(1)}, & \cdots, & x_D^{(1)} \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)}, & x_2^{(N)}, & \cdots, & x_D^{(N)} \end{bmatrix} \begin{bmatrix} w^{(0)} \\ w^{(1)} \\ w^{(2)} \\ \vdots \\ w^{(D)} \end{bmatrix}$$

Matrix form

instead of $\hat{y}^{(n)} = \mathbf{w}^T \mathbf{x}^{(n)}$
 $\in \mathbb{R}$ $1 \times D$ $D \times 1$

use **design matrix** to write $\hat{\mathbf{y}} = \mathbf{X} \mathbf{w}$
 $N \times 1$ $N \times D$ $D \times 1$



```
yh = np.dot(X, w)  
cost = np.sum((yh - y)**2)
```

Linear least squares

$$J(w) = \frac{1}{2} \|y - Xw\|^2 = \frac{1}{2} (y - Xw)^T (y - Xw)$$

squared L2 norm of the **residual** vector

Normal equation

system of D linear equations

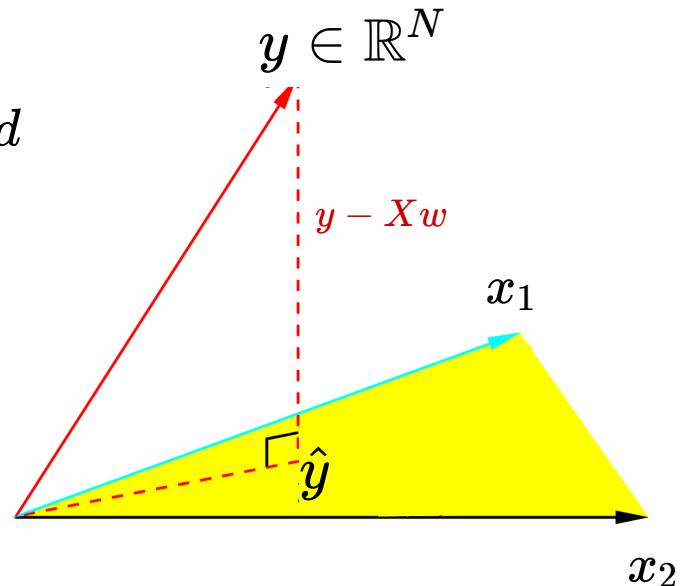
$$\sum_n (y^{(n)} - w^T x^{(n)}) x_d^{(n)} = 0 \quad \forall d$$

matrix form (using the design matrix)

each row enforces one of D equations

$D \times N$ $N \times 1$

$$X^T (y - Xw) = \vec{0}$$



Normal equation: because for optimal w , the residual vector is
normal to column space of the design matrix

Direct solution

we can get a closed form solution!

$$X^T(y - Xw) = \vec{0}$$

$$X^T X w = X^T y$$

pseudo-inverse of X

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

$D \times D$ $D \times N$ $N \times 1$

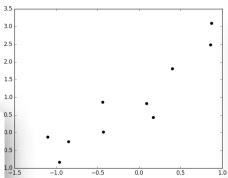
$$\hat{y} = Xw = X(X^T X)^{-1} X^T y$$

projection matrix into column space of X

Example

```
1 X=array([[ 0.86],      1 Y=array([[ 2.49],  
2 [ 0.09],      2 [ 0.83],  
3 [-0.85],      3 [-0.25],  
4 [ 0.87],      4 [ 3.1 ],  
5 [-0.44],      5 [ 0.87],  
6 [-0.43],      6 [ 0.02],  
7 [-1.1 ],      7 [-0.12],  
8 [ 0.4 ],      8 [ 1.81],  
9 [-0.96],      9 [-0.83],  
10 [ 0.17]]) ) 10 [ 0.43]])  
11  
12  
13 X = np.hstack((np.ones((X.shape[0],1)),X))  
14  
15 array([[ 1. ,  0.86],  
16 [ 1. ,  0.09],  
17 [ 1. , -0.85],  
18 [ 1. ,  0.87],  
19 [ 1. , -0.44],  
20 [ 1. , -0.43],  
21 [ 1. , -1.1 ],  
22 [ 1. ,  0.4 ],  
23 [ 1. , -0.96],  
24 [ 1. ,  0.17]])
```

Todo: try this at home.



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

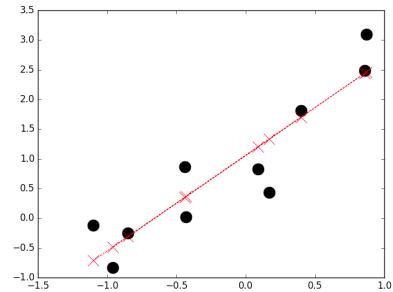
```
1 X.T.dot(X)  
2 >> array([[ 10.      , -1.39     ],  
3 [ -1.39     ,  4.9261   ]])  
4 np.linalg.inv(X.T.dot(X))  
5 >> array([[ 0.10408228,  0.02936895],  
6 [ 0.02936895,  0.2112874  ]])  
7  
8 w = np.linalg.inv(X.T.dot(X)).dot(X.T.dot(Y))  
9 array([[ 1.05881341],  
10 [ 1.61016842]])  
11
```

```
w = np.linalg.lstsq(X,Y)[0]
```



```
1 yh = X.dot(w)  
2 >> array([[ 2.44355825],  
3 [ 1.20372857],  
4 [ -0.30982974],  
5 [ 2.45965993],  
6 [ 0.35033931],  
7 [ 0.36644099],  
8 [ -0.71237185],  
9 [ 1.70288078],  
10 [ -0.48694827],  
11 [ 1.33254204]])
```

$$y = 1.06 + 1.61x$$



Time complexity

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

$\mathcal{O}(D^3)$ matrix inversion
 $\mathcal{O}(D^2 N)$ $D \times D$ elements, each requiring N multiplications

$\mathcal{O}(ND)$ D elements, each using N ops.

total complexity for $N > D$ is $\mathcal{O}(D^3 + ND^2)$

total complexity for $N > D$ is $\mathcal{O}(ND^2)$

in practice we don't directly use matrix inversion (unstable)

Multiple targets

instead of $y \in \mathbb{R}^N$ we have $Y \in \mathbb{R}^{N \times D'}$

a different weight vectors for each target

each column of Y is associated with a column of W

$$\hat{Y} = XW$$

$N \times D$ $N \times D$ $D \times D'$

$$W^* = (X^T X)^{-1} X^T Y$$

$D \times D$ $D \times N$ $N \times D'$



```
w = np.linalg.lstsq(X,Y)[0]
```

Nonlinear basis functions

So far we learned a linear function $f_w = \sum_d w_d x_d$

nothing changes if we have nonlinear bases $f_w = \sum_d w_d \phi_d(x)$

solution simply becomes $w^* = (\Phi^T \Phi)^{-1} \Phi^T y$

replacing X with Φ

a (nonlinear) feature

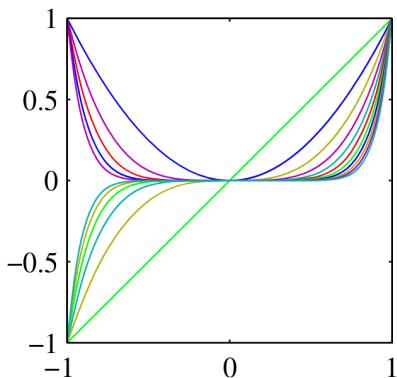
$$\Phi = \begin{bmatrix} \phi_1(x^{(1)}), & \phi_2(x^{(1)}), & \cdots, & \phi_D(x^{(1)}) \\ \phi_1(x^{(2)}), & \phi_2(x^{(2)}), & \cdots, & \phi_D(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(N)}), & \phi_2(x^{(N)}), & \cdots, & \phi_D(x^{(N)}) \end{bmatrix}$$

one instance

Nonlinear basis functions

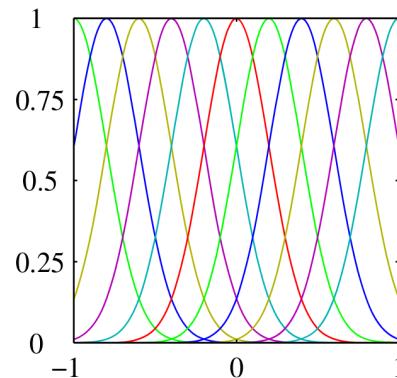
examples

original input is scalar $x \in \mathbb{R}$



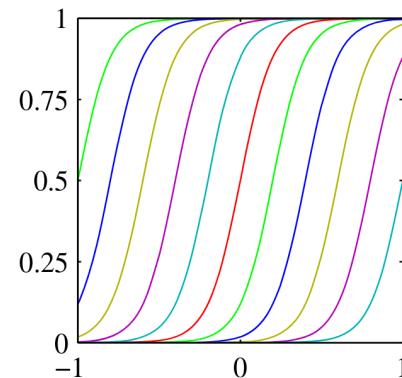
polynomial bases

$$\phi_k(x) = x^k$$



Gaussian bases

$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$



Sigmoid bases

$$\phi_k(x) = \frac{1}{1+e^{-\frac{x-\mu_k}{s}}}$$

Nonlinear basis: example

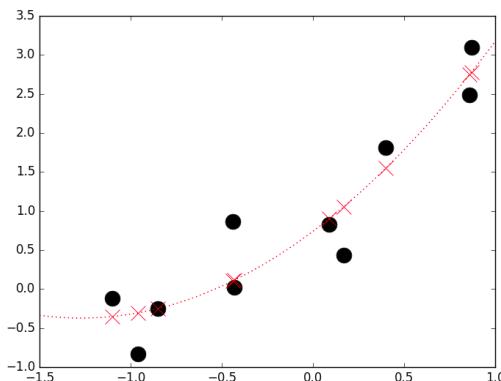
```
1 X=array([[ 1. ,  0.86],      1 Y=array([[ 2.49],
2   [ 1. ,  0.09],      2   [ 0.83],
3   [ 1. , -0.85],      3   [-0.25],
4   [ 1. ,  0.87],      4   [ 3.1 ],
5   [ 1. , -0.44],      5   [ 0.87],
6   [ 1. , -0.43],      6   [ 0.02],
7   [ 1. , -1.1 ],      7   [-0.12],
8   [ 1. ,  0.4 ],      8   [ 1.81],
9   [ 1. , -0.96],      9   [-0.83],
10  [ 1. ,  0.17]]))       10  [ 0.43]]))
```

```
w = np.linalg.lstsq(X,Y)[0]  
array([[ 0.73960888],  
       [ 1.74736998],  
       [ 0.68670054]])
```

$$y = 0.74 + 1.75x + 0.68x^2$$

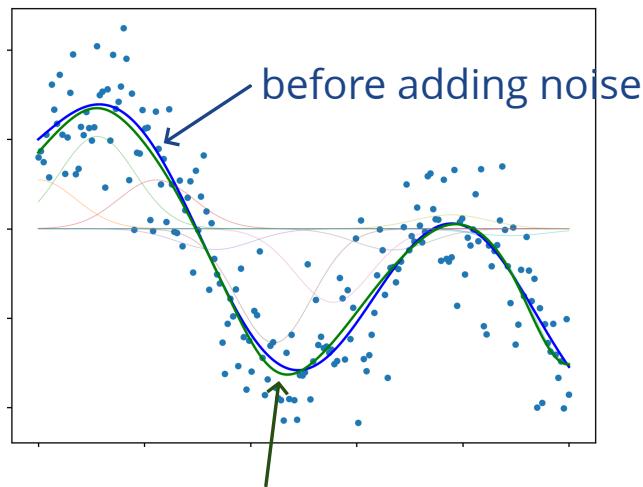
```
1 yh = x.dot(w)  
2 >> array([[ 2.75023077],  
3   [ 0.90243445],  
4   [-0.24951447],  
5   [ 2.77958439],  
6   [ 0.10371131],  
7   [ 0.11521071],  
8   [-0.35159045],  
9   [ 1.54842895],  
10  [-0.30500309],  
11  [ 1.05650742]])
```

Todo: try this at home and add x^3, x^4, \dots .



Gaussian bases example

$$\phi_k(x) = e^{-(x-\mu_k)^2}$$



$$y^{(n)} = \sin(x^{(n)}) + \cos(\sqrt{|x^{(n)}|}) + \epsilon$$

...

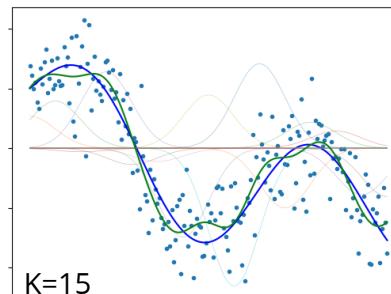
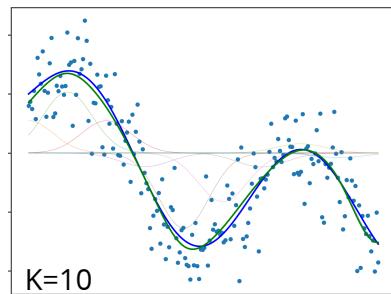
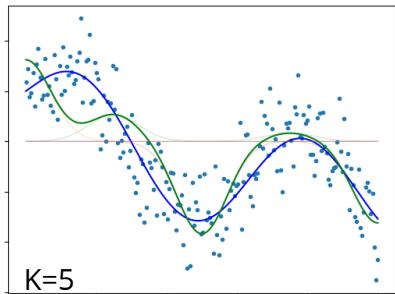
```
1 #x: N
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-(x-mu)**2)
5 mu = np.linspace(0,10,10) #10 Gaussians bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```

Todo: try this at home

Gaussian bases example

$$\phi_k(x) = e^{-(x-\mu_k)^2}$$

$$y^{(n)} = \sin(x^{(n)}) + \cos(\sqrt{|x^{(n)}|}) + \epsilon$$



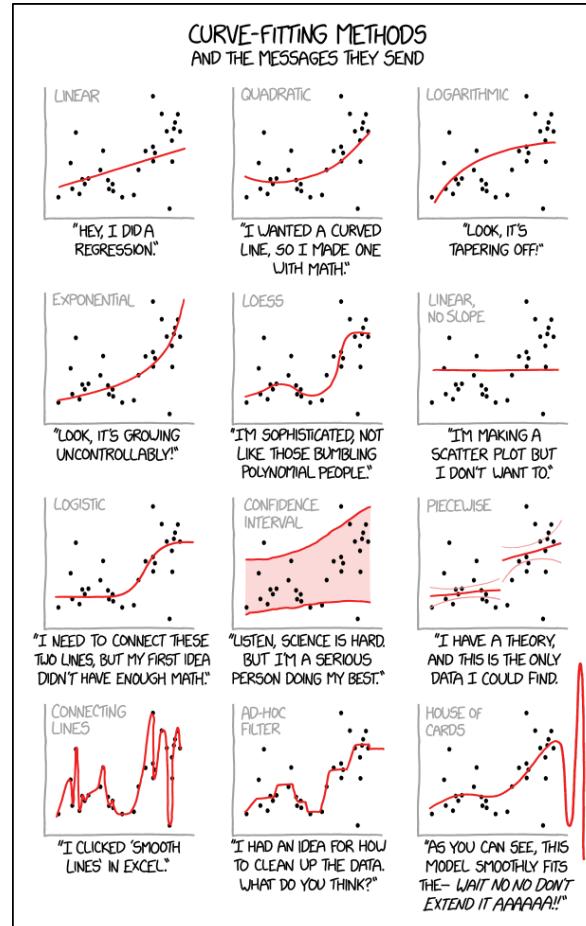
Todo: try this at home, change number
of bases, see how the fit changes.



```
1 #x: N
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-(x-mu)**2)
5 mu = np.linspace(0,10,K) #K Gaussians bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```

For fun!

https://www.explainxkcd.com/wiki/index.php/2048:_Curve-Fitting



Problematic settings

$$\text{In } W^* = (X^T X)^{-1} X^T Y$$

what if we have a **large dataset**?

- use stochastic gradient descent

what if $X^T X$ is **not invertible**?

columns of X are not linearly independent
(either redundant features or D > N)

- W^* is **not unique, make it unique** by
 - removing redundant features
 - regularization (*later!*)



```
w = np.linalg.lstsq(X,Y)[0]
```

- **or find one** of the solutions
 - decomposition-based (not discussed) methods still work
 - use gradient descent (*later!*)

Summary

linear regression:

- models targets as a linear function of features
- fit the model by minimizing sum of squared errors
- has a direct solution with $\mathcal{O}(D^3 + ND^2)$ complexity
- gradient descent: future

we can build complex models:

- using any number of non-linear features
- ensure features are linearly independent