

# REGRESSION PROBLEMS

① In a regression problem, we have a set of evidence variables  $X[1], \dots, X[D]$  and a response variable  $Y$  that we want to predict.

For instance, let's say we want to predict cholesterol level given age and weight:

often called  
the "bias"

X (evidence vars)

X[1]  
(offset)

X[2]  
(age)

X[3]  
(weight)

$X_1 = [$

1

24

150

$]$

$X_2 = [$

1

50

164

$]$

$\vdots$

$\vdots$

$X_N = [$

1

22

205

$]$

Y (response var)

(cholesterol)

182

=

$Y_1$

210

=

$Y_2$

$\vdots$

$\vdots$

202

=

$Y_N$

We have  $N$  training examples, each consisting of a vector  $X_i$  and a scalar  $Y_i$ , to learn from.

## REGRESSION PROBLEMS

2) We assume the response variable is generated using the following steps:

- we start with a vector of weights (one weight per evidence variable), e.g.

$$W = \begin{bmatrix} W[1] \\ W[2] \\ W[3] \end{bmatrix} = \begin{bmatrix} -50 \\ 2 \\ 1 \end{bmatrix}$$

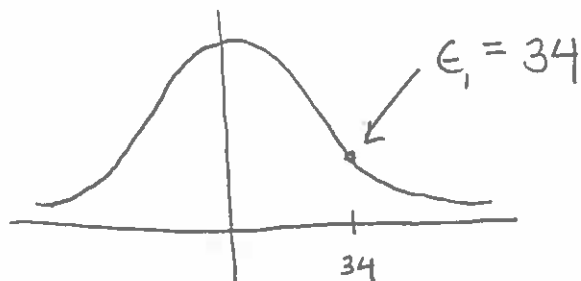
- we compute the weighted linear combination of the evidence variables (this will result in a single number):

$$W^T X_1 = [W[1] \ W[2] \ W[3]] \begin{bmatrix} X_1[1] \\ X_1[2] \\ X_1[3] \end{bmatrix}$$

the offset allows us to shift the total up or down by a constant factor

$$\begin{aligned} &= -50 \cdot 1 + 2 \cdot 24 + 1 \cdot 150 \\ &= 148 \end{aligned}$$

- we sample a random number  $\epsilon$  from a distribution  $\psi$  e.g.  $\epsilon \sim \text{Normal}(0, \sigma^2)$

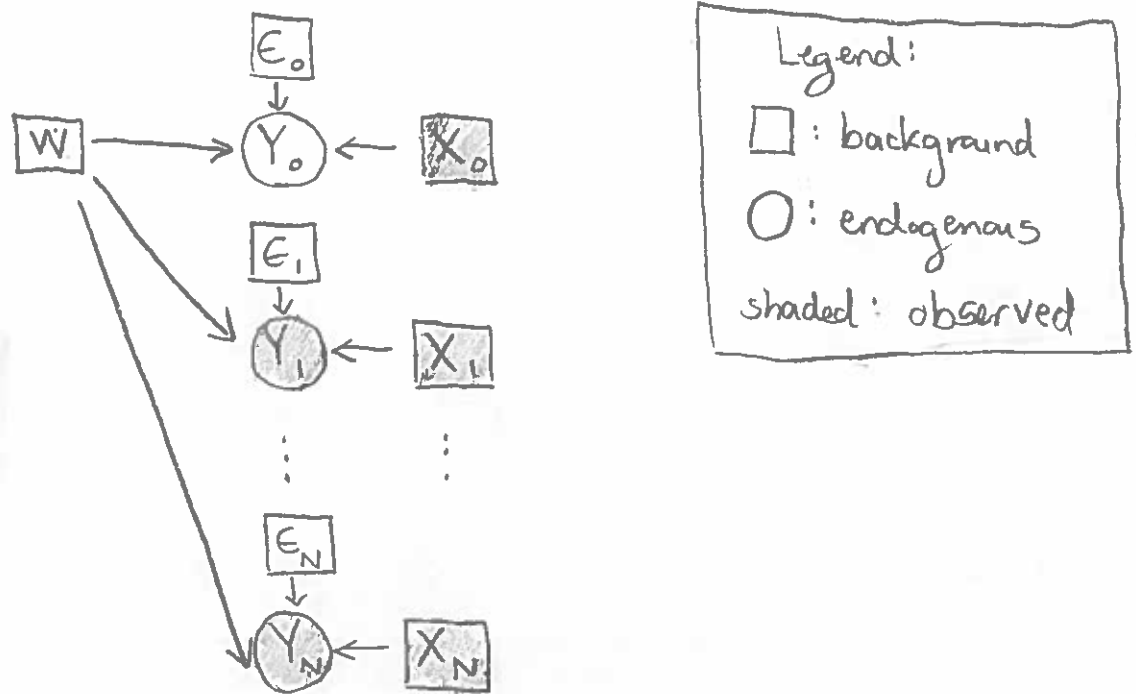


- we compute the response variable as a function  $p$  of  $W^T X$  and  $\epsilon$ , e.g.

$$Y = p(W^T X_1, \epsilon_1) = W^T X_1 + \epsilon_1 = 148 + 34 = 182$$

## REGRESSION PROBLEMS

- ③ Given this setup, we want to predict the value of an unobserved response variable  $Y_0$  given its observed evidence variables  $X_0$  and our training data. In its simplest formulation, the causal diagram looks as follows:



We assume a probability distribution  $P$  over the background variables such that all background variables are marginally independent, i.e.

$$P(w, \epsilon_0, \dots, \epsilon_n, x_0, \dots, x_n) = P_w(w) P_\epsilon(\epsilon_0) \dots P_\epsilon(\epsilon_n) \cdot P_x(x_0) \dots P_x(x_n)$$

Moreover we assume that all variables  $\epsilon_n$  are drawn from the same distribution  $P_\epsilon$  and all variables  $x$  are drawn from the same distribution  $P_x$ .

## REGRESSION PROBLEMS

④ What's the deal with the  $\epsilon_i$ 's?

We'll call these stochastic terms. The idea is to allow some softness around the deterministic point  $w^T x_i$ . In other words, just because  $w^T x_i = 148$  for age=24 and weight=150, that doesn't mean we want the model to claim that EVERY person whose age is 24 and whose weight is 150 must have a cholesterol level of EXACTLY 148.

Instead, we want their cholesterol levels to be dispersed around 148.

⑤ Suppose we choose  $P_\epsilon \sim \text{Normal}(0, \sigma^2)$  for some fixed variance  $\sigma^2$ , and suppose we choose  $\rho(z, \epsilon) = z + \epsilon$ . This is what we did in ②, which gave us a cholesterol level of  $148 + 34 = 182$  for a 24-year-old, 150lb subject.

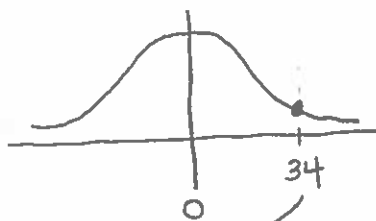
DETERMINISTIC TERM

$$W^T X = \begin{bmatrix} -50 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 24 \\ 150 \end{bmatrix}$$

$$= 148$$

STOCHASTIC TERM

$$\epsilon =$$



$$\oplus$$

$$182$$

## REGRESSION PROBLEMS

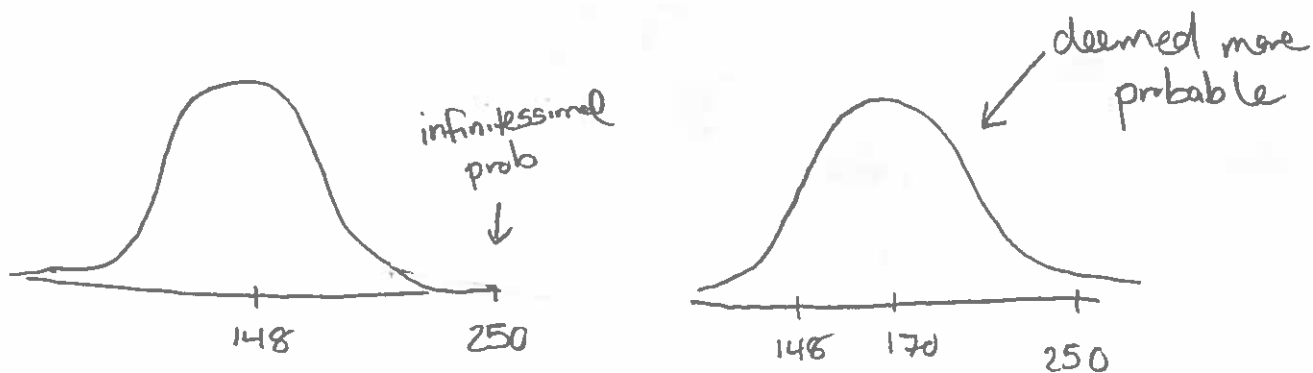
⑥ These choices:

$$\begin{cases} P_{\epsilon} \sim \text{Normal}(0, \sigma^2) \\ \rho(z, \epsilon) = z + \epsilon \end{cases}$$

give us the ordinary linear regression model

⑦ But it may not be the best choice. One possible downside of the normal distribution is that it has rapidly diminishing tails, so a normal distribution centered at 148 may give a nearly infinitesimal probability to 250.

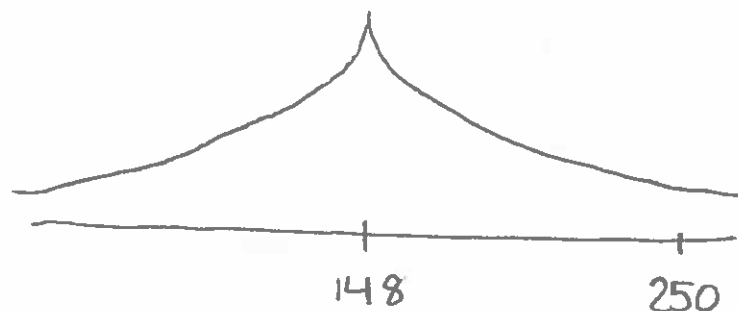
In practice, this means that using the normal distribution makes regression very sensitive to outliers and noise, because even you have a thousand 24-year-olds with cholesterol 148, just one 24-year-old with cholesterol 250 can cause the learned mean to shift significantly upwards, because this will be deemed more probable than even one person with cholesterol 250:



## REGRESSION PROBLEMS

---

- ⑧ If this is a problem for you, you can use a "heavy-tailed" distribution (i.e. they diminish in probability much more slowly). One example is the Laplace distribution.



Often, however, such distributions are not as computationally convenient.

- 
- ⑨ This choice:

$$\begin{cases} P_{\epsilon} \sim \text{Laplace}(0, b) \\ \rho(z, \epsilon) = z + \epsilon \end{cases}$$

gives us the "robust" linear regression model

## REGRESSION PROBLEMS

- ⑩ Another common choice for  $P_e$  and  $p$  comes into play when the response variable is Boolean-valued, e.g.:

X (evidence vars)			Y (response var)				
X[1]	X[2]	X[3]					
(offset)	(age)	(weight)	(has high cholesterol)				
$x_1 = [$	1	24	150	$]$	0	=	$y_1$
$x_2 = [$	1	50	164	$]$	1	=	$y_2$
$\vdots$					$\vdots$		$\vdots$
$x_N = [$	1	22	205	$]$	1	=	$y_N$

- ⑪ We could use ordinary linear regression, but then we end up predicting values in the range  $(-\infty, \infty)$ , rather than restricting ourselves to the set  $\{0, 1\}$ .

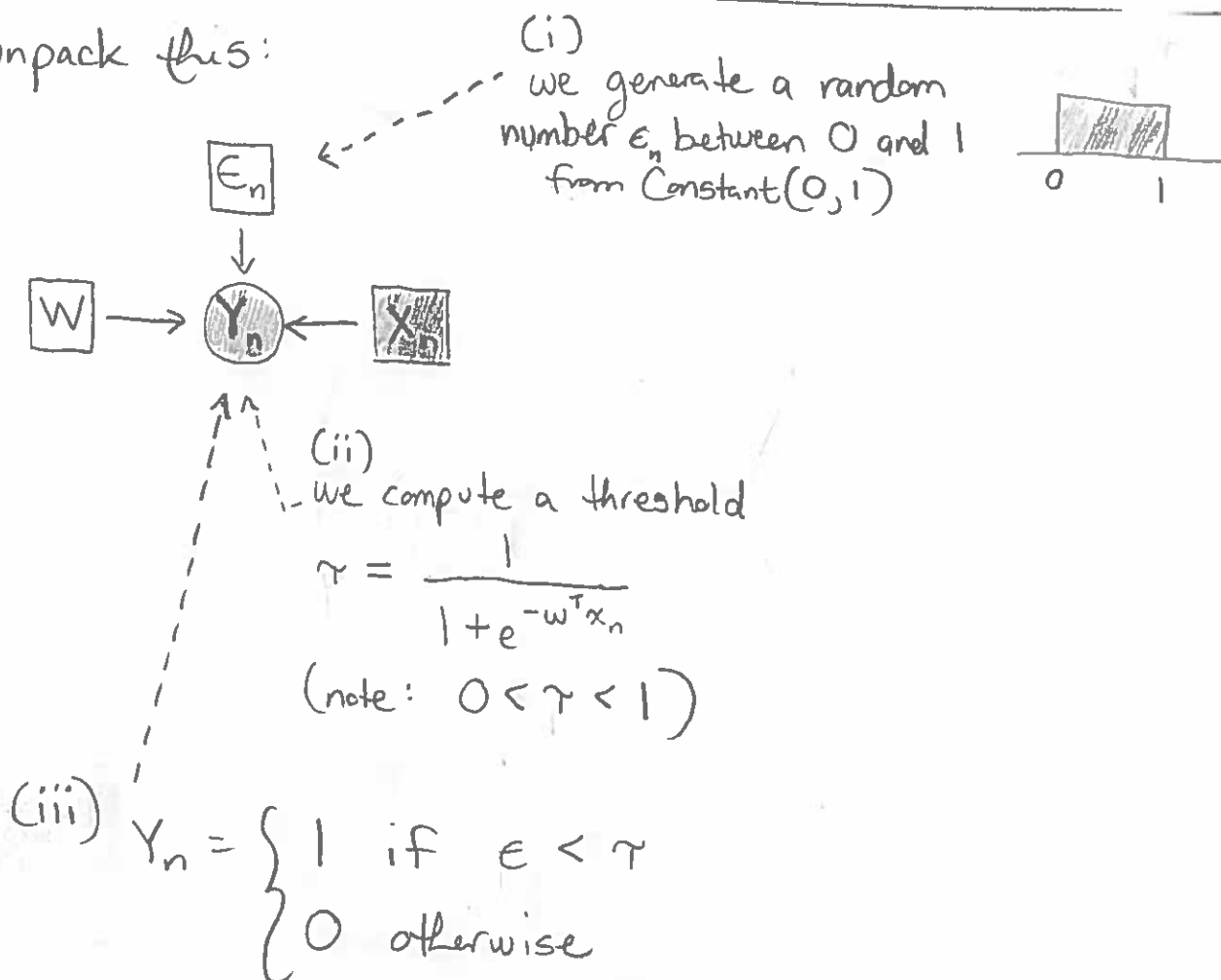
Another choice:

$$\begin{cases} P_e \sim \text{Constant}(0, 1) \\ p(z, e) = \frac{1}{(1 + e^{-z})^{-1} < e(z)} \end{cases}$$

gives us the famous logistic regression model

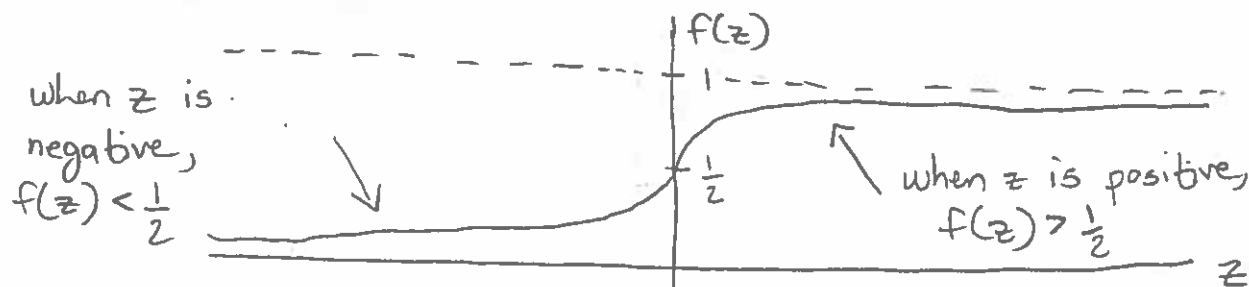
## REGRESSION PROBLEMS

⑫ Let's unpack this:



In other words, the probability that  $Y_n = 1$  is equal to  $\frac{1}{1 + e^{-w^T x_n}}$ .

⑬ This function,  $f(z) = (1 + e^{-z})^{-1}$ , looks like this:



is called the logistic function (or the logit, or the sigmoid).



## REGRESSION PROBLEMS

---

- ⑭ No matter what regression model we're using, we typically want to predict the value of the unobserved response variable  $Y_0$  given its observed evidence variables  $X_0$  (recall  $X_0$  is a vector) and our other observations (i.e.  $X_n, Y_n$  for  $n \geq 1$ ).

In other words, we want  $\operatorname{argmax}_{y_0} P(y_0 | x_0, x_1, \dots, x_N, y_1, \dots, y_N)$ .

---

- ⑮ The exact computation would be:

$$\begin{aligned} & \operatorname{argmax}_{y_0} P(y_0 | x_0, x_1, \dots, x_N, y_1, \dots, y_N) \\ &= \operatorname{argmax}_{y_0} \int P(y_0, w | x_0, x_1, \dots, x_N, y_1, \dots, y_N) dw \end{aligned}$$

[Law of Total Probability]

But integrals are painful to work with, so let's not go there.

## REGRESSION PROBLEMS

⑩ Instead, we can do a point estimate approach.

(a) compute the most probable value of  $\hat{w}$  given the observations:

$$\hat{w} = \operatorname{argmax}_w P(w | x_0, x_1, \dots, x_N, y_1, \dots, y_N)$$

$$= \operatorname{argmax}_w \frac{P(x_0, x_1, \dots, x_N, y_1, \dots, y_N | w) P(w)}{P(x_0, x_1, \dots, x_N, y_1, \dots, y_N)} \quad [\text{Bayes Rule}]$$

$$= \operatorname{argmax}_w P(x_0, x_1, \dots, x_N, y_1, \dots, y_N | w) P(w) \quad [\text{remove constant factors}]$$

$$= \operatorname{argmax}_w P(x_0 | w) P(x_1 | x_0, w) \dots P(y_N | x_0, \dots, x_N, y_1, \dots, y_{N-1}, w) P(w) \quad [\text{Chain Rule of Prob}]$$

$$= \operatorname{argmax}_w P(x_0) P(x_1) \dots P(x_N) P(y_1 | w, x_1) \dots P(y_N | w, x_N) P(w) \quad [\text{d-separation}]$$

$$= \operatorname{argmax}_w P(y_1 | w, x_1) \dots P(y_N | w, x_N) P(w) \quad [\text{remove constant factors}]$$

$$= \operatorname{argmax}_w P(w) \prod_{n=1}^N P(y_n | w, x_n)$$

(b) compute the most probable value of  $\hat{y}_0$  given  $\hat{w}$ :

$$\hat{y}_0 = \operatorname{argmax}_{y_0} P(y_0 | \hat{w}, x_0, x_1, \dots, x_N, y_1, \dots, y_N)$$

$$= \operatorname{argmax}_{y_0} P(y_0 | \hat{w}, x_0) \quad [\text{d-separation}]$$

# REGRESSION PROBLEMS

⑦ In short:

(a) compute  $\hat{w} = \operatorname{argmax}_w P(w) \prod_{n=1}^N P(y_n | w, x_n)$

(b) compute  $\hat{y}_0 = \operatorname{argmax}_{y_0} P(y_0 | \hat{w}, x_0)$

This is called the MAP (maximum a posteriori) estimate.

⑧ A special case of the MAP estimate assumes that  $P(w)$  is the same for every possible  $w$ . Since it then becomes a constant factor, we can drop it from the  $\operatorname{argmax}$ :

(a) compute  $\hat{w} = \operatorname{argmax}_w \prod_{n=1}^N P(y_n | w, x_n)$

(b) compute  $\hat{y}_0 = \operatorname{argmax}_{y_0} P(y_0 | \hat{w}, x_0)$

This is called the MLE (maximum likelihood estimate).

how can this be, if there's an infinite space of  $w$ -vectors? don't worry about it

