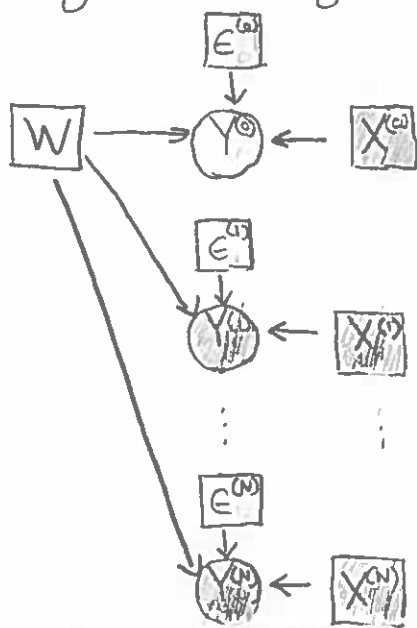


LINEAR REGRESSION: MAP

① Recall "ordinary linear regression":



where: $P_{\epsilon}(\epsilon^{(n)}) \sim \text{Normal}(0, \sigma^2) \quad \forall n \in \{0, \dots, N\}$

$$y^{(n)} \leftarrow w^T x^{(n)} + \epsilon^{(n)}$$

② Also recall that one way to estimate the value of the unobserved response variable $Y^{(0)}$ is through maximum a posteriori (MAP) estimation:

(a) compute $\hat{w} = \underset{w}{\operatorname{argmax}} P(w) \prod_{n=1}^N P(y^{(n)} | w, x^{(n)})$

(b) compute $\hat{y}^{(0)} = \underset{y^{(0)}}{\operatorname{argmax}} P(y^{(0)} | \hat{w}, x^{(0)})$

In the MLE approach, we assume that all weight vectors are equally likely (without further evidence), so we treat $P(w)$ as a constant and drop it from the equation.

LINEAR REGRESSION: MAP

- ③ But maybe we do have an opinion about which weight vectors are more likely prior to observing any evidence (this is called a prior probability or an a priori belief).

First off, why would we have such an opinion?

- ④ Consider if we actually wanted to predict someone's cholesterol accurately on the basis of lifestyle factors. We don't know what might be relevant, so we throw a lot of evidence variables into the mix:

X (evidence vars)						Y (response var)
X_1	X_2	X_3	X_4	...	X_{10000}	
(offset)	(age)	(weight)	(smoking freq)		(gumchewing freq)	(cholesterol)

- ⑤ Most of these evidence vars probably don't have any impact on cholesterol, so we expect that a good weight vector $w = \begin{bmatrix} w_1 \\ \vdots \\ w_{10000} \end{bmatrix}$ will contain mostly

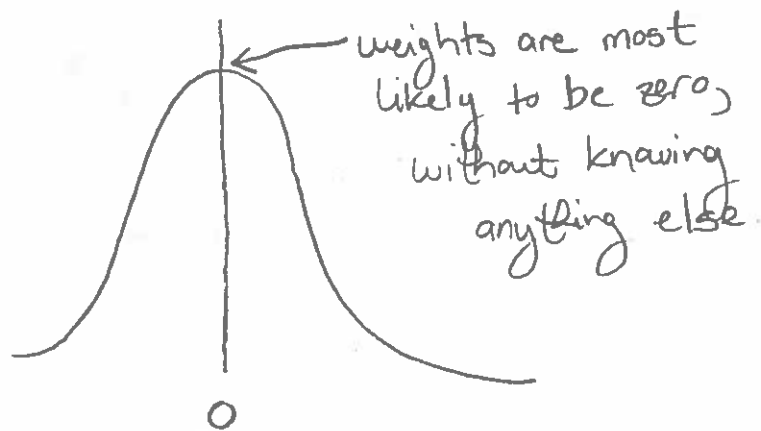
zeros, since then $Y = w^T x$ will only be a function of a small subset of the evidence vars.

LINEAR REGRESSION: MAP

- ⑥ So our apriori belief is that weights are most likely to be zero.

How can we express this as a distribution?

One way is to say that for each weight w_d ,
 $P(w_d) \sim \text{Normal}(0, \sigma^2)$ for some variance σ^2 :



-
- ⑦ So let's see if we can simplify ②(a), our point estimate:

$$\begin{aligned}\hat{w} &= \underset{w}{\operatorname{argmax}} P(w) \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \underset{w}{\operatorname{argmax}} \log P(w) \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \underset{w}{\operatorname{argmax}} \ell(w)\end{aligned}$$

LINEAR REGRESSION: MAP

⑧ Continuing:

$$\begin{aligned} l(w) &= \log P(w) + \sum_{n=1}^N \log P(y^{(n)} | w, x^{(n)}) \\ &= \log P(w) + l_{MLE}(w) \end{aligned}$$

where $l_{MLE}(w)$ is the likelihood function for the MLE (see LINEAR REGRESSION: MLE, ④)

⑨ As with ordinary linear regression, we'll assume the stochastic terms $\epsilon^{(n)}$ are normally distributed, i.e.

$$P_{\epsilon} \sim \text{Normal}(0, \sigma^2)$$

We'll also use the prior distribution over weights that we argued for in ⑥:

$$P(w) \sim \text{Normal}(0, \tau^2)$$

not necessarily the same variance



These choices give us a type of regression called ridge regression.

LINEAR REGRESSION: MAP

⑩ Continuing to simplify with these choices:

$$\hat{w} = \operatorname{argmax}_w \ell(w)$$

$$= \operatorname{argmax}_w \log P(w) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \log \left(\prod_{d=1}^D \left(\frac{1}{2\pi\tau^2} \right)^{\frac{1}{2}} \exp \left(\frac{-1}{2\tau^2} w_d^2 \right) \right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \sum_{d=1}^D \log \left[\left(\frac{1}{2\pi\tau^2} \right)^{\frac{1}{2}} \exp \left(\frac{-1}{2\tau^2} w_d^2 \right) \right] + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \left(\sum_{d=1}^D -\frac{1}{2} \log 2\pi\tau^2 \right) + \left(\sum_{d=1}^D -\frac{1}{2\tau^2} w_d^2 \right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \left[-\frac{D}{2} \log 2\pi\tau^2 - \frac{1}{2\tau^2} \sum_{d=1}^D w_d^2 \right] + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \frac{-\sum_{d=1}^D w_d^2}{2\tau^2} + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \frac{-w^T w}{2\tau^2} + \ell_{MLE}(w)$$

LINEAR REGRESSION: MAP

11) From LINEAR REGRESSION: MLE, (5), we know that:

$$\begin{aligned} \ell_{\text{MLE}}(w) &= -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2 \\ &= -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) \end{aligned}$$

Plugging this in to what we have so far:

$$\hat{w} = \underset{w}{\operatorname{argmax}} \quad \frac{-w^T w}{2\tau^2} - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)$$

$$= \underset{w}{\operatorname{argmax}} \quad -\frac{1}{2\tau^2} w^T w - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)$$

$$= \underset{w}{\operatorname{argmax}} \quad -\frac{\sigma^2}{\tau^2} w^T w - (y - Xw)^T (y - Xw)$$

$$= \dots \quad [\text{see LINEAR REGRESSION: MLE (8)}]$$

$$= \underset{w}{\operatorname{argmax}} \quad -\frac{\sigma^2}{\tau^2} w^T w + 2w^T X^T y - w^T X^T X w$$

$$= \underset{w}{\operatorname{argmax}} \quad 2w^T X^T y - w^T X^T X w - \frac{\sigma^2}{\tau^2} w^T w$$

$$= \underset{w}{\operatorname{argmin}} \quad w^T X^T X w - 2w^T X^T y + \frac{\sigma^2}{\tau^2} w^T w$$

LINEAR REGRESSION: MAP

⑫ So the loss function for ridge regression is:

$$L_{\text{ridge}}(w) = w^T X^T X w - 2w^T X^T y + \frac{\sigma^2}{\tau^2} w^T w$$

Notice that this can be expressed in terms of the loss function for ordinary linear regression

$$L_{\text{lin}}(w) = w^T X^T X w - 2w^T X^T y$$

$$L_{\text{ridge}}(w) = L_{\text{lin}}(w) + \frac{\sigma^2}{\tau^2} w^T w$$

⑬ That means ridge regression's loss function is simply the usual linear regression loss function, plus some constant multiple of the squared " L_2 -norm" of the weight vector:

$$\hat{w} = \underset{w}{\operatorname{argmin}} L_{\text{ridge}}(w)$$

$$= \underset{w}{\operatorname{argmin}} \underbrace{L_{\text{lin}}(w)} + K \cdot \underbrace{w^T w}$$

we want the likelihood of the data to be high

but we also want the weights to be close to zero

$$w^T w = \sum_{d=1}^D w_d^2 = \|w\|_2^2$$



LINEAR REGRESSION: MAP

⑭ Thus ridge regression's loss function is combining two different objectives:

(a) we want the likelihood of the training data to be high

$$\operatorname{argmin}_w L_{\text{lin}}(w)$$

(b) we want the learned weight vector to have a small L_2 -norm (i.e. we want the length of the weight vector to be small)

$$\operatorname{argmin}_w \|w\|_2^2$$

Objective (b) is often called a regularization term and so ridge regression is sometimes known as linear regression with L_2 -regularization.

⑮ So going back to ②, ridge regression estimates the value of the unobserved response variable $y^{(i)}$ as follows:

$$\begin{aligned} \text{(a) compute point estimate } \hat{w} &= \operatorname{argmax}_w P(w) \prod_{n=1}^N P(y^{(n)} | w, x^{(n)}) \\ &= \operatorname{argmin}_w L_{\text{ridge}}(w) \end{aligned}$$

$$\text{(b) compute } \hat{y}^{(i)} = \operatorname{argmax}_{y^{(i)}} P(y^{(i)} | \hat{w}, x^{(i)})$$

$$= \hat{w}^T x^{(i)} \quad (\text{see LINEAR REGRESSION: MLE, ③})$$

LINEAR REGRESSION: MAP

⑩ Exercise: Adapt LINEAR REGRESSION: MLE ⑩-⑪ to compute a closed-form expression for $\arg\min_w L_{\text{ridge}}(w)$.

$$\begin{aligned}\text{Solution: } \nabla L_{\text{ridge}}(w) &= \frac{\partial}{\partial w} \left(L_{\text{lin}}(w) + \frac{\sigma^2}{\tau^2} w^T w \right) \\&= \frac{\partial}{\partial w} L_{\text{lin}}(w) + \frac{\sigma^2}{\tau^2} \frac{\partial}{\partial w} w^T w \\&= \frac{\partial}{\partial w} L_{\text{lin}}(w) + \frac{\sigma^2}{\tau^2} \begin{bmatrix} \frac{\partial}{\partial w_1} w^T w \\ \vdots \\ \frac{\partial}{\partial w_D} w^T w \end{bmatrix} \\&= \frac{\partial}{\partial w} L_{\text{lin}}(w) + \frac{\sigma^2}{\tau^2} \begin{bmatrix} 2w_1 \\ \vdots \\ 2w_D \end{bmatrix} \\&= \frac{\partial}{\partial w} L_{\text{lin}}(w) + \frac{2\sigma^2}{\tau^2} w \\&= 2X^T X w - 2X^T y + \frac{2\sigma^2}{\tau^2} w\end{aligned}$$

$$\text{Then } \nabla L_{\text{ridge}}(w) = 0$$

$$\Rightarrow 2X^T X w - 2X^T y + \frac{2\sigma^2}{\tau^2} w = 0$$

$$\Rightarrow \left(X^T X + \frac{\sigma^2}{\tau^2} I \right) w = X^T y$$

$$\Rightarrow w = \left(X^T X + \frac{\sigma^2}{\tau^2} I \right)^{-1} X^T y$$

LINEAR REGRESSION: MAP

⑦ This gives us the following algorithm for ridge regression:

RIDGE REGRESSION (X, y , ratio $\frac{\sigma^2}{\tau^2}$, input $x^{(0)}$):

- compute point estimate $\hat{w} = (X^T X + \frac{\sigma^2}{\tau^2} I)^{-1} X^T y$

- compute estimate $\hat{y}^{(0)} = \hat{w}^T x^{(0)}$

- return $\hat{y}^{(0)}$

⑧ That's what we get if we assume the weights are drawn from a $\text{Normal}(0, \tau^2)$ distribution.

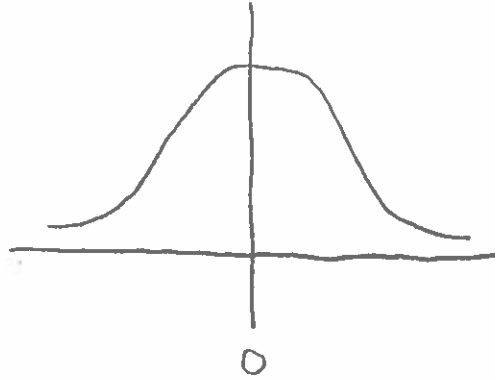
But the whole motivation was that we had a bunch of evidence variables, most of which were completely irrelevant:

X (evidence vars)					Y (response var)
X_1	X_2	X_3	X_4	X_{10000}	
(offset)	(age)	(weight)	(smoking freq)	(gumchewing freq)	(cholesterol)

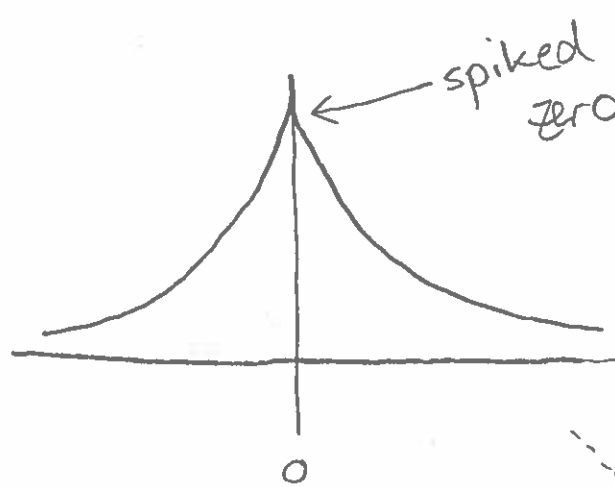
Ideally, we'd want most of the weights to be exactly zero, not just close to zero. That way, the nonzero weights tell us which features are relevant.

LINEAR REGRESSION: MAP

- ①9 So rather than having a prior distribution that softly favors zero:



we'd like a prior distribution that strongly favors zero:



it's Laplace!



-
- ②0 So let's assume:

$$P_{\epsilon} \sim \text{Normal}(0, \sigma^2)$$

$$P(w_d) \sim \text{Normal}(0, \tau^2) \text{ Laplace}(0, b)$$

These choices give us a type of regression called lasso.

LINEAR REGRESSION: MAP

21) We computed the point estimate \hat{w} for ridge regression in (10). If we use Laplace instead, we get:

$$\hat{w} = \operatorname{argmax}_w \ell(w)$$

$$= \operatorname{argmax}_w \log P(w) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \log \prod_{d=1}^D \frac{1}{2b} \exp\left(-\frac{|w_d|}{b}\right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \sum_{d=1}^D \log\left(\frac{1}{2b} \exp\left(-\frac{|w_d|}{b}\right)\right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \left(\sum_{d=1}^D \log \frac{1}{2b} + \log \exp\left(-\frac{|w_d|}{b}\right) \right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w \left(\sum_{d=1}^D -\log 2b \right) - \left(\sum_{d=1}^D \frac{|w_d|}{b} \right) + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w -D \log 2b - \frac{1}{b} \sum_{d=1}^D |w_d| + \ell_{MLE}(w)$$

$$= \operatorname{argmax}_w -\frac{1}{b} \sum_{d=1}^D |w_d| + \ell_{MLE}(w)$$

LINEAR REGRESSION: MAP

22) From LINEAR REGRESSION: MLE, (5), we know that:

$$l_{MLE}(w) = -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)$$

Plugging this in to what we have so far:

$$\hat{w} = \operatorname{argmax}_w -\frac{1}{b} \sum_{d=1}^D |w_d| - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)$$

$$= \operatorname{argmax}_w -\frac{1}{b} \sum_{d=1}^D |w_d| - \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)$$

$$= \operatorname{argmax}_w \left(-\frac{2\sigma^2}{b} \sum_{d=1}^D |w_d| \right) - (y - Xw)^T (y - Xw)$$

$$= \dots \text{ [see LINEAR REGRESSION: MLE (7)]}$$

$$= \operatorname{argmax}_w \left(-\frac{2\sigma^2}{b} \sum_{d=1}^D |w_d| \right) + 2w^T X^T y - w^T X^T X w$$

$$= \operatorname{argmax}_w 2w^T X^T y - w^T X^T X w - \frac{2\sigma^2}{b} \sum_{d=1}^D |w_d|$$

$$= \operatorname{argmin}_w w^T X^T X w - 2w^T X^T y + \frac{2\sigma^2}{b} \sum_{d=1}^D |w_d|$$

LINEAR REGRESSION: MAP

② So the loss function for lasso regression is:

$$L_{\text{lasso}}(w) = w^T X^T X w - 2w^T X^T y + \frac{2\sigma^2}{b} \sum_{d=1}^D |w_d|$$
$$= L_{\text{lin}}(w) + \frac{2\sigma^2}{b} \sum_{d=1}^D |w_d|$$

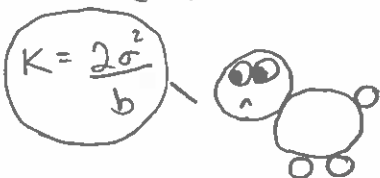
The expression $\sum_{d=1}^D |w_d|$ is called the L_1 -norm of weight vector w :

$$\|w\|_1 = \sum_{d=1}^D |w_d|$$

Thus lasso's loss function is simply the usual linear regression loss function, plus some constant multiple of the L_1 -norm of the weight vector:

$$\hat{w} = \underset{w}{\operatorname{argmin}} L_{\text{lasso}}(w)$$

$$= \underset{w}{\operatorname{argmin}} L_{\text{lin}}(w) + K \cdot \|w\|_1$$



we want the likelihood
of the data to be high

but we also want the
weights to be close
to zero

Lasso is sometimes called linear regression with L_1 -regularization.

LINEAR REGRESSION: MAP

②④ In summary, lasso regression estimates the value of an unobserved response variable Y_0 as follows:

(a) compute point estimate $\hat{w} = \underset{w}{\operatorname{argmin}} L_{\text{lasso}}(w)$

(b) compute $\hat{y}_0 = \hat{w}^T x_0$

②⑤ To compute $\underset{w}{\operatorname{argmin}} L_{\text{lasso}}(w)$, compute:

$$\nabla L_{\text{lasso}}(w) = \frac{\partial}{\partial w} \left(L_{\text{lin}}(w) + K \sum_{d=1}^D |w_d| \right)$$

$$= \frac{\partial}{\partial w} L_{\text{lin}}(w) + K \frac{\partial}{\partial w} \sum_{d=1}^D |w_d|$$

$$= \frac{\partial}{\partial w} L_{\text{lin}}(w) + K \begin{bmatrix} \frac{\partial}{\partial w_1} \sum_{d=1}^D |w_d| \\ \vdots \\ \frac{\partial}{\partial w_D} \sum_{d=1}^D |w_d| \end{bmatrix}$$

$$= \frac{\partial}{\partial w} L_{\text{lin}}(w) + K \begin{bmatrix} \operatorname{sgn}(w_1) \\ \vdots \\ \operatorname{sgn}(w_D) \end{bmatrix}$$

$$= 2X^T X w - 2X^T y + K \begin{bmatrix} \operatorname{sgn}(w_1) \\ \vdots \\ \operatorname{sgn}(w_D) \end{bmatrix}$$

LINEAR REGRESSION: MAP

- ② Now we could try to find the critical points of $L_{\text{lasso}}(w)$ by setting $\nabla L_{\text{lasso}}(w) = 0$ and solving for w :

$$2X^T X w - 2X^T y + K \begin{bmatrix} \text{sgn}(w_1) \\ \vdots \\ \text{sgn}(w_D) \end{bmatrix} = 0$$

$w = ???$

but ick.

- ⑦ Fortunately, we have like 6 variants of gradient descent to choose from, since we can compute $\nabla L_{\text{lasso}}(w)$.

Thus we have an algorithm for lasso:

LASSO REGRESSION

- compute point estimate $\hat{w} = \text{GRADDESCENT}(L_{\text{lasso}})$
- compute estimate $\hat{y}^{(0)} = \hat{w}^T x^{(0)}$
- return $\hat{y}^{(0)}$