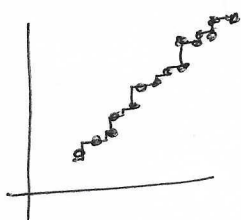# Linear Regression

Regression → Predicting a continuous variable given some other variables.



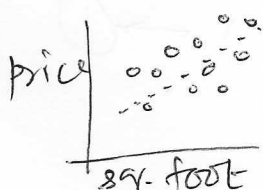← basically an interpolation. (doesn't work well in most of the cases).

## Linear fitting to data.

We want to fit a linear function to an observed set of points $X = [x_1, \cdots x_N]$ with associated labels $Y = [y_1, \cdots, y_N]$. Once we fit the function, we want to use it to predict the $y$ for new $X$.



price | sq. foot

search for line that best fits these data points.

I/P: $\vec{x} \in \mathbb{R}^d$      The model we'll use:

O/P: $y \in \mathbb{R}$

$$\hat{y} = W_0 + W_1 X_1 + \cdots + W_d X_d$$

$$= [1 \ X_1 \cdots X_d] \begin{bmatrix} W_0 \\ \vdots \\ W_d \end{bmatrix}$$

$$\underset{bias}{\nearrow} = X^T W \quad (\text{o } W^T X \to \text{why ?})$$

why?

However, we may not find the true mapping function $f: X \to Y$, but we'll try to approximate it as much as possible.

## Residuals/errors : $e_i = y_i - \hat{y}_i$

A natural loss function is squared loss:

$$L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2 \quad \begin{array}{l} \text{\& we want this to be as} \\ \text{small as possible} \end{array}$$

## Least squares estimation : Minimize w.r.t. $w$

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

Many ways to solve this (e.g. taking partial derivatives $\frac{\partial L}{\partial W_0} = 0$, $\frac{\partial L}{\partial W_1} = 0, \cdots$ & solve a system of eqn's.). We'll use matrix notation & avoid solving system of eqn's.

Let,

$$X = \begin{bmatrix} - & \vec{x_1} & - \\ - & \vec{x_2} & - \\ & \vdots & \\ - & \vec{x_n} & - \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ & \vdots & & \\ 1 & x_{N1} & \cdots & x_{Nd} \end{bmatrix} \qquad \begin{array}{l} X \Rightarrow N \times (d+1) \\ Y \Rightarrow N \times 1 \\ W \Rightarrow (d+1) \times 1. \end{array}$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad W = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}. \quad \hat{Y} = XW$$

Hence,

$$L(W) = \tfrac{1}{N} \|Y - \hat{Y}\|_2^2 \quad (L_2 \text{ norm squarred})$$

$$= \tfrac{1}{N} \|Y - XW\|_2^2$$

$$= \tfrac{1}{N}(Y - XW)^T(Y - XW)$$

$$L(W) = \tfrac{1}{N}\left[ Y^TY + W^TX^TXW - 2Y^TXW \right] \quad (\text{how?!}) \longrightarrow \text{①}$$

Now we need to take $\dfrac{\partial L(W)}{\partial W} = 0$

But $\vec{W}$ is a vector, so we need to take care of the derivatives.
We'll get back to eq$^n$ ①. Let's see some derivatives first.

$$\frac{\partial(\vec{W}^T\vec{X})}{\partial\vec{W}} = \left[ \frac{\partial(\vec{W}^T\vec{X})}{\partial W_0} \cdots \frac{\partial(\vec{W}^T\vec{X})}{\partial W_d} \right]$$

$$\frac{\partial\left[\sum_{j=0}^{d} W_j x_j\right]}{\partial W_0} = x_0 \quad (=1 \text{ in this case})$$

so, $\dfrac{\partial(\vec{W}^T\vec{X})}{\partial\vec{W}} = [x_0 \ x_1 \cdots x_d] \Rightarrow \boxed{\dfrac{\partial(\vec{W}^T\vec{X})}{\partial\vec{W}} = \vec{X}^T}$

or, $\boxed{\dfrac{\partial(X^TW)}{\partial W} = X^T} \longrightarrow \text{②}$

Now let's compute $\dfrac{\partial}{\partial\vec{W}}\left(\vec{W}^TA\vec{W}\right)$, where $A$ is a symmetric matrix

$$\frac{\partial(\vec{W}^TA\vec{W})}{\partial\vec{W}} = \left[ \frac{\partial(W^TAW)}{\partial W_0} \cdots \frac{\partial(W^TAW)}{\partial W_d} \right]$$

$$\frac{\partial}{\partial W_0}(W^TAW) = \frac{\partial}{\partial W_0}\left[ \sum_{i=0}^{d}\sum_{j=0}^{d} W_i a_{ij} W_j \right]$$

$$\Rightarrow \frac{\partial}{\partial w_0}\left(\vec{w}^T A \vec{w}\right) = \frac{\partial}{\partial w_0}\left[a_{00} w_0^2 + \sum_{i \neq 0} a_{i0} w_i w_0 + \sum_{j \neq 0} a_{0j} w_0 w_j \right.$$
$$\left. + \sum_{i \neq 0}\sum_{j \neq 0} w_i a_{ij} w_j \right]$$

$$= 2 a_{00} w_0 + \sum_{i \neq 0} a_{i0} w_i + \sum_{j \neq 0} a_{0j} w_j$$

$$= 2 a_{00} w_0 + 2 \sum_{i \neq 0} a_{i0} w_i \qquad (\text{how?!})$$

$$= 2 \sum_{i=0}^{d} a_{i0} w_i$$

$$= 2 \vec{w}^T A_0 \qquad \left(A_0 \text{ is the first column of } A \to A_0 = \begin{pmatrix} a_{00} \\ a_{10} \\ \vdots \\ a_{d0} \end{pmatrix}\right)$$

$$\Rightarrow \frac{\partial(\vec{w}^T A \vec{w})}{\partial \vec{w}} = \left[2\vec{w}^T A_0 \quad 2\vec{w}^T A_1 \quad \cdots \quad 2\vec{w}^T A_d\right]$$
$$= \boxed{2\vec{w}^T A} \longrightarrow ③$$

Now let's get back to eqⁿ ①:

$$\frac{\partial L(w)}{\partial \vec{w}} = \frac{\partial}{\partial \vec{w}}\left[\frac{1}{n}\left(\vec{y}^T \vec{y} - 2\vec{y}^T x \vec{w} + \vec{w}^T x^T x \vec{w}\right)\right]$$
$$= \frac{1}{n}\left(-2\vec{y}^T x + 2\vec{w}^T x^T x\right) \qquad (\text{using } 2, 3)$$
$$= \frac{1}{n}\left(-2\vec{y}^T x + 2\vec{w}^T x^T x\right) = 0$$

Now set $\frac{\partial L}{\partial \vec{w}} = 0 \Rightarrow \frac{1}{n}\left(-2\vec{y}^T x + 2\vec{w}^T x^T x\right) = 0$

$$\Rightarrow \vec{w}^T x^T x = \vec{y}^T x$$
$$\Rightarrow (x^T x) \vec{w} = x^T \vec{y} \qquad (\text{taking transpose})$$
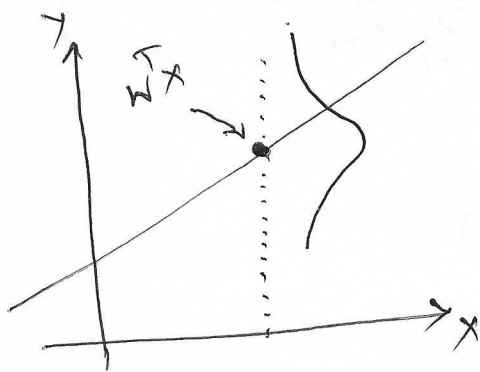
$$\Rightarrow \boxed{\hat{w}_{OLS} = (x^T x)^{-1} x^T \vec{y}}$$

$(x^T x)^{-1} x^T$ is called the pseudo-inverse of $x$.
$\downarrow$

$$\det x^+ = (x^T x)^{-1} x^T$$
$$\text{Then } x^+ x = (x^T x)^{-1} x^T x = I$$

# Probabilistic view of linear regression.



Let $X \sim P(X)$ (some unknown distribution)

& let, $y|X=x \sim \mathcal{N}(w^T x, \sigma^2)$

The error is,

$$e = y - w^T x$$
$$\sim \mathcal{N}(0, \sigma^2) \quad (\text{zero-mean Gaussian})$$

Let's see the MLE estimator of $W$ (assume $\sigma^2$ is constant)

$$\hat{W}_{MLE} = \arg\max_{\vec{w}} \log P(D/\vec{w})$$

$$= \arg\max_{\vec{w}} \sum_{i=1}^{N} \log P(e_i/\vec{w})$$

$$= \arg\max_{\vec{w}} \sum_{i=1}^{N} \left[ -\frac{(y_i - w^T x_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

$\uparrow$ const.

$\downarrow$ doesn't depend on $W$

Dataset:

$$\{(x_1, y_1) \cdots \}$$
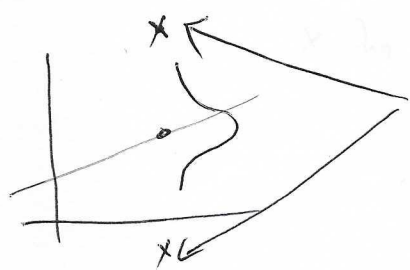
convert into errors

$e_1 \cdots e_n$

since errors are Gaussians,

$$P(e) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[ -\frac{(y - w^T x)^2}{2\sigma^2} \right]$$

$$\Rightarrow \boxed{\begin{array}{c} \hat{W}_{MLE} = \arg\min_{\vec{w}} \sum_{i=1}^{N} (y_i - w^T x_i)^2 \\[2mm] \hat{W}_{MLE} = \hat{W}_{OLS} \end{array}}$$
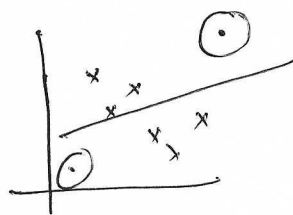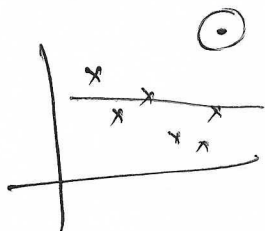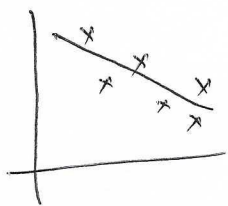
So the assumption we're making is that errors are coming from a standard Gaussian.



probability of points to be here is very small. If there are points like this, linear regression fails.

OLS is not robust to these points, aka outliers.

Our model doesn't consider "outliers", so it
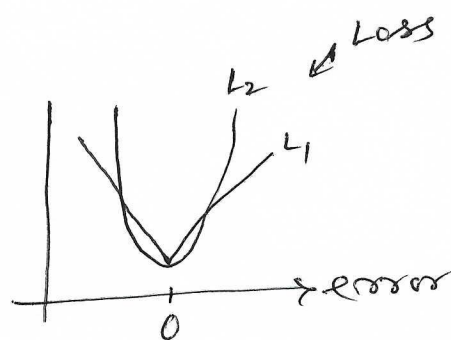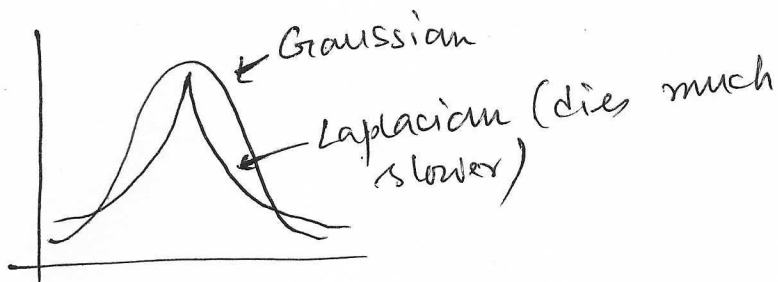performs poorly.



↓ how to fix these?

Instead of Gaussian,

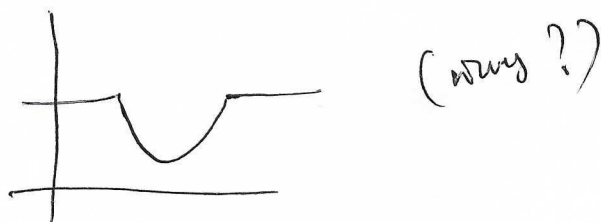$$Y | X = x = \text{Laplacian} (W^T x, b)$$

$$Y \sim \text{Laplacian} (\mu, b)$$

$$\Rightarrow P(y | \mu, b) = \frac{1}{2b} \exp \left[ - \frac{|y - \mu|}{b} \right]$$

$$\hat{W}_{MLE} = \arg \min_W \sum |y_i - W^T x_i|$$



← Gaussian

← Laplacian (dies much slower)



$L_2$   ↙ Loss

$L_1$

→ error

0

Q A robust loss function should look like:



(why?)

Why don't we do that?
  ↳ Convexity!! extremely hard to optimize.