Start simple w/ <u>**Linear Regression**</u>

Let's continue using house price example but use an input or feature in $\mathbb{R}^2$ - so sq ft and number of bedrooms.

So $x$'s $\in \mathbb{R}^2$ and for notation

$$X_1^{(i)} = \text{sq ft of house } (i)$$

$$X_2^{(i)} = \text{\# of bedrooms in house } (i)$$

Now, want to decide what should look like.

Here, we are doing linear regression (really affine) in X.

i.e. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x_2$

where $\theta_i$ are the parameters in our models (sometimes called weights)

So regression here means choosing the "best" weights

Note $h_\theta(x) : X \rightarrow Y$

$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

often will write $h$, not $h_\theta$, sometimes even $h(x; \theta)$.

Depends on if we want to emphasize parameter dependence.

Also, for linear regression sometimes introduce

$x_0 = 1$ (intercept form)

$d$ ← # of input variables, not counting intercept one

$$h(x) = \sum_{i=0}^{d} \theta_i x_i = \theta^T x = \theta \cdot x = (\theta, x)$$

So how do we choose the $\theta_i$?

Will do dumb but generalizable way, then linear algebra way.

The $\theta_i$'s depend on how we define a good fit.

Have some sense $h(x)$ should be close to $y$ for our training set — seems reasonable.

Thus, we write down a measure of goodness of fit.

Call it:
- Loss function
- Cost            Etc.
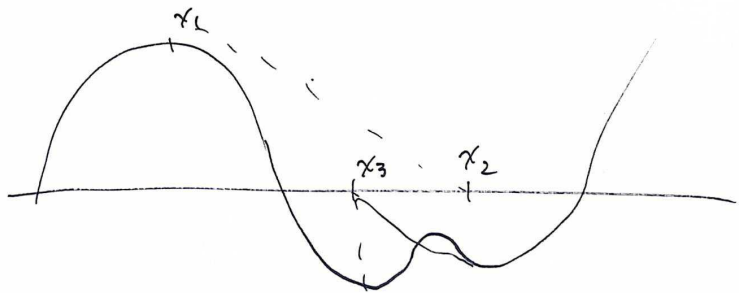- Energy
- Error

## Arbitrary Choice:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

And the goal is to find $\theta$ to make this as small as possible.

(This special choice of $J$ called $\boxed{\text{ordinary least squares)}}$

So want $\theta$ to minimize $J$

or $\theta$ that satisfies $\underset{\theta}{\text{argmin}} \, J(\theta) = \theta$

$\qquad\qquad\qquad\qquad$ Argument That Minimizes.

Want to describe an algorithm to find $\theta$:

Start w/ initial guess for $\theta$ and use algorithm that changes $\theta$. Slowly hopefully making $J$ smaller at each step.
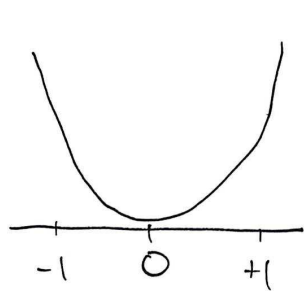
# Gradient Descent

Start w/ $\theta$ guess and want formula for next $\tilde{\theta}$.

$$\tilde{\theta}_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

This is update for $j^{th}$ component but would update $\theta_j$ simultaneously

Basically, $---$ is a gradient: $\tilde{\theta} = \tilde{\theta} - \alpha \nabla J(\theta)$

Interestingly, $\alpha$ is called the learning rate.



$x^2$       $\nabla J(x) = 2x$

$-\alpha 2$

$1 - 1 \cdot 2$

The gradient of $J$ gives the direction in which $J(\theta)$ is increasing the fastest, so $-\nabla J(\theta)$ is steepest descent.

Thus, we need to figure out $\nabla J(\theta)$.

# Gradient Descent for Least Squares

$$\tilde{\Theta} = \Theta - \alpha \nabla J(\Theta)$$

Need to figure out $\nabla J(\Theta)$

Recall $J(\Theta) = \frac{1}{2} \sum_{i=1}^{n} \left( h_\Theta(x^{(i)}) - y^{(i)} \right)^2$

$$h_\Theta(x) = \sum_{i=0}^{d} \Theta_i x_i = \Theta^T x$$

$x_0 = 1$

Start w/ $n=1$, so no sum in $T$

training set is $(x, y)$

$$\frac{\partial}{\partial \Theta_j} J(\Theta) = \frac{\partial}{\partial \Theta_j} \frac{1}{2} \left( h_\Theta(x) - y \right)^2$$

$$= 2 \cdot \frac{1}{2} \left( h_\Theta(x) - y \right) \frac{\partial}{\partial \Theta_j} \left( h_\Theta(x) - y \right)$$

$$= \left( h_\Theta(x) - y \right) \frac{\partial}{\partial \Theta_j} \left( \sum_{i=0}^{d} \Theta_i x_i - y \right)$$

$$= \left( h_\Theta(x) - y \right) x_j$$

If only have one training example,
then the update is

$$\tilde{\Theta}_j = \Theta_j + \alpha \left( y^{(i)} - h_\Theta(x^{(i)}) \right) x_j^{(i)}$$

Note: magnitude of update $\underset{\underset{\text{proportional to}}{\nearrow}}{\alpha}$ error term $\left( y^{(i)} - h_\Theta(x^{(i)}) \right)$

If your prediction nearly predicts value of $y^{(i)}$,
then little update.

What about more data? "$\sum$" is linear or is
the derivative

$$\tilde{\Theta}_j = \Theta_j + \alpha \sum_{i=1}^{n} \left( y^{(i)} - h_\Theta(x^{(i)}) \right) x_j^{(i)} \quad , \quad \forall j$$

OR

$$\underset{\sim}{\Theta} = \underset{\sim}{\Theta} + \alpha \sum_{i=1}^{n} \left( y^{(i)} - h_\Theta(x^{(i)}) \right) x^{(i)}$$

scalars

vectors

called batch gradient descent
- uses all training data in every step
  for update to the parameters $\underset{\sim}{\Theta}$.

$$h_\Theta(x): \underset{\mathbb{R}^n}{x} \to \underset{\mathbb{R}}{y}$$

J is in this case convex so should have
a unique solution ; gradient descent
should converge.

Can modify gradient in various ways to get good
optimization schemes.

As an alternative, we can use the single
update rule but w/ more data.

for $i = 1, \ldots, n$

$$\tilde{\Theta} = \Theta + \alpha\left(y^{(i)} - h_\Theta(x^{(i)})\right) x^{(i)}$$

Repeat until whole thing has converged.

So we just keep updating $\tilde{\Theta}$ one data point
at a time – called

Stochastic Gradient Descent
(incremental)

Update as data is introduced instead of using
all the data.