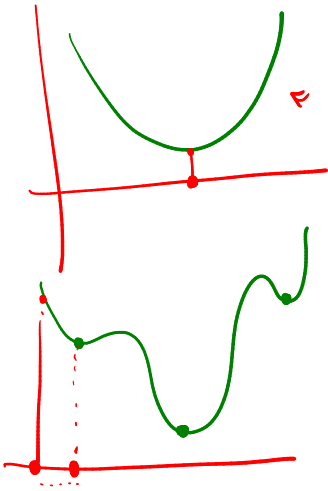


COMPUTATIONAL CONSIDERATIONS

$$w^* = (x x^T)^+ x y$$

inverse computation
is expensive if d is large
 $O(d^3)$



- We know w^* is the solution of an unconstrained optimization

- We can apply GRADIENT DESCENT.

$$w^{t+1} = w^t - \eta^t \nabla f(w^t)$$

η^t \rightarrow scalar step size
 $\nabla f(w^t)$ \rightarrow gradient at w^t

$$f(w) = \|xw - y\|^2 = \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\nabla f(w) = 2(x x^T)w - 2xy \quad [\text{verify this}]$$

Gradient descent update for Linear regression

$$w^{t+1} = w^t - \eta^t [2(x x^T)w - 2xy]$$

- What if n is large [hundreds of millions], might want to avoid $x x^T$
- How to adapt gradient descent?

STOCHASTIC GRADIENT DESCENT

for $t=1, \dots, T$

- At each step, sample a bunch of datapoints (\tilde{x}) uniformly at random from the set of all points
- PRETEND this sample is the entire dataset and take a gradient step w.r.t it

$$2(\tilde{x}\tilde{x}^T w^t - \tilde{x}\tilde{y})$$

↳ Manageable because $\tilde{x} \in \mathbb{R}^{d \times 1}$ \hookrightarrow #sampled points.

end

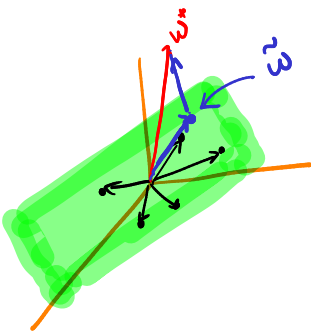
After T rounds, use

$$\underline{w_{\text{sgd}}^T} = \frac{1}{T} \sum_{i=1}^T w^t$$

as opposed to w^T as in Standard Gradient descent

Guaranteed to converge to optima with high probability

NON-LINEAR REGRESSION



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

- $\underline{w^*}$ must lie in the span of data points.

$$\sum_{i=1}^n (\underline{w^*}^T x_i - y_i)^2 = \sum_{i=1}^n (\underline{\tilde{w}}^T x_i - y_i)^2$$

$$\begin{aligned} w^* &= \tilde{w} + \tilde{w}_\perp \\ \text{+i} \quad \underline{w^*}^T x_i &= (\underline{\tilde{w}} + \underline{\tilde{w}_\perp})^T x_i = \tilde{w}^T x_i + \underbrace{\tilde{w}_\perp^T x_i}_0 \end{aligned}$$

$$\omega^* = \boxed{X \alpha^*} \quad \text{for some } \alpha^* \in \mathbb{R}^n$$

$$= (X X^T)^T X y$$

$$X \alpha^* = (X X^T)^T X y$$

$$(X X^T) X \alpha^* = (X X^T) (X X^T)^T X y$$

$$(X X^T) X \alpha^* = X y$$

$$\underbrace{X^T (X X^T)}_{X^T X} X \alpha^* = \underbrace{X^T X}_{X^T X} y$$

$$\underbrace{(X^T X)^2}_{:= K} \alpha^* = (X^T X) y$$

$$\Rightarrow K^2 \alpha^* = K y$$

$$\boxed{\alpha^* = K^{-1} y} \leftarrow \leftarrow$$

$\in \mathbb{R}^{n \times n}$

KERNEL - REGRESSION

PREDICTION

for some $x_{\text{test}} \in \mathbb{R}^d$

$$\omega^{*T} \phi(x_{\text{test}})$$

$$= \left(\sum_{i=1}^n \alpha_i^* \phi(x_i) \right)^T \phi(x_{\text{test}})$$

$$= \sum_{i=1}^n \alpha_i^* \underbrace{\phi(x_i)^T \phi(x_{\text{test}})}_{K(x_i, x_{\text{test}})}$$

How important
is input
towards
 ω^*

How
similar
is x_{test}
to x_i

PROBABILISTIC VIEW OF
LINEAR REGRESSION

$$x \in \mathbb{R}^d \quad y \in \mathbb{R}$$

Dataset

$$\{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$y/x \sim \boxed{W}x + \epsilon$$

Unknown
but fixed
 $\epsilon \in \mathbb{R}^d$

Noise
 $N(0, \sigma^2)$
↳ Gaussian.

- Can view this as an "ESTIMATION" problem
- Solution approach - Maximum Likelihood.

Likelihood

$$L(\mathbf{w}; \mathbf{x}_1, \dots, \mathbf{x}_n, y_1, \dots, y_n) = \prod_{i=1}^n e^{-\frac{(\mathbf{w}^T \mathbf{x}_i - y_i)^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

$$\log L(\mathbf{w}; \mathbf{x}_1, \dots, \mathbf{x}_n, y_1, \dots, y_n) = \sum_{i=1}^n -\frac{(\mathbf{w}^T \mathbf{x}_i - y_i)^2}{2\sigma^2} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

equivalently

$$\max_{\mathbf{w}} \sum_{i=1}^n -(\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

$$= \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

$$\hat{\mathbf{w}}_{ML} = \mathbf{w}^* = (\mathbf{X}\mathbf{X}^T)^{\dagger} \mathbf{X}\mathbf{y}$$

CONCLUSION: Maximum Likelihood estimator assuming ZERO MEAN GAUSSIAN NOISE is same as LINEAR REGRESSION with SQUARED ERROR!

What else have we gained?

- Can study properties of estimators $\hat{\mathbf{w}}_{ML}$!