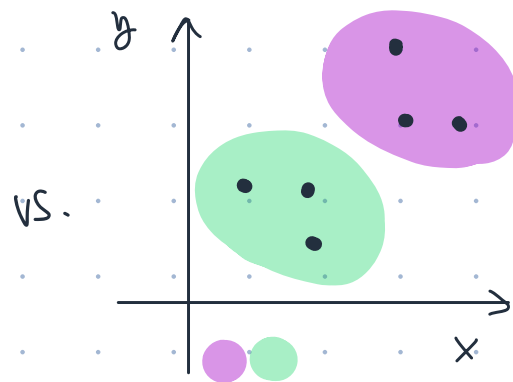
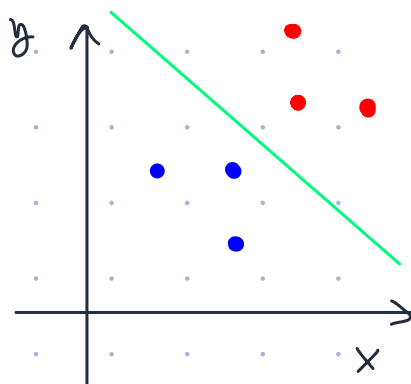


ANNOUNCEMENTS

1. HW3 due Friday, ^{11:59 pm} late due Sunday - NO extensions beyond late due
 2. Prelim conflict form - fill by 3:30 pm today!
 3. P3 due changed from 03/04 → 03/05 11:59 pm
 4. Prelim logistics - will also be posted to Ed later!
 5. HW4 to be released w/ solutions (soon!)
-

TIME TO TURN YOUR NON-NOTE-TALKING DEVICES OFF!

SO FAR

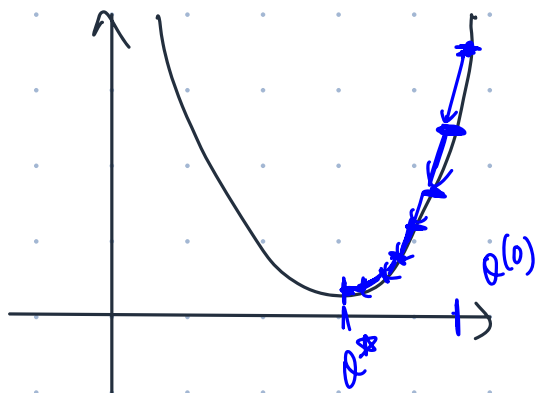


REGIME

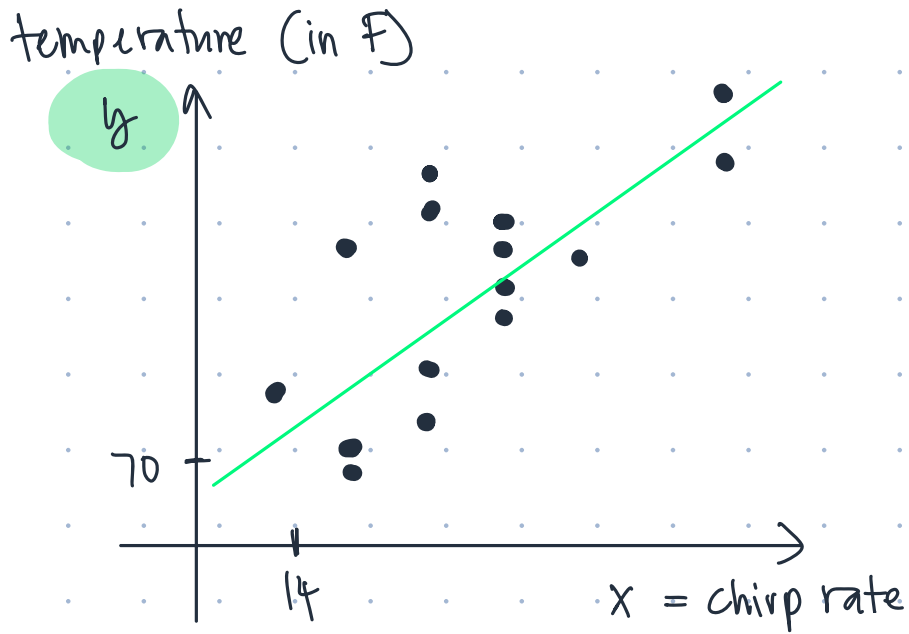
Supervised learning — has labels
Unsupervised learning — no labels

Estimating "cost" function — MLE, MAP

OPTIMIZATION — Iterative approaches: GD, Adagrad, momentum, Newton



TODAY : linear regression — instead of a discrete class, we wish to predict a continuous "y"



$$\text{temperature} = \theta_0 + \theta_1 \text{ chirp rate} + \epsilon \quad \text{— unmodelled noise}$$

GOAL : Can we recover the "green" line?

$$\text{temperature} = \theta_0 + \theta_1 \times \text{chirp rate}$$

we are trying to learn some "h"

$$\begin{aligned} h(x^{(j)}; \theta) &= \theta_0 \cdot 1 + \theta_1 x_1^{(j)} + \theta_2 x_2^{(j)} + \dots + \theta_d x_d^{(j)} \\ &= \sum_{i=0}^d \theta_i x_i^{(j)} = \theta^T \tilde{x}^{(j)}, \quad \tilde{x}^{(j)} = \begin{bmatrix} x_0^{(j)} = 1 \\ x_1^{(j)} \\ \vdots \\ x_d^{(j)} \end{bmatrix} \end{aligned}$$

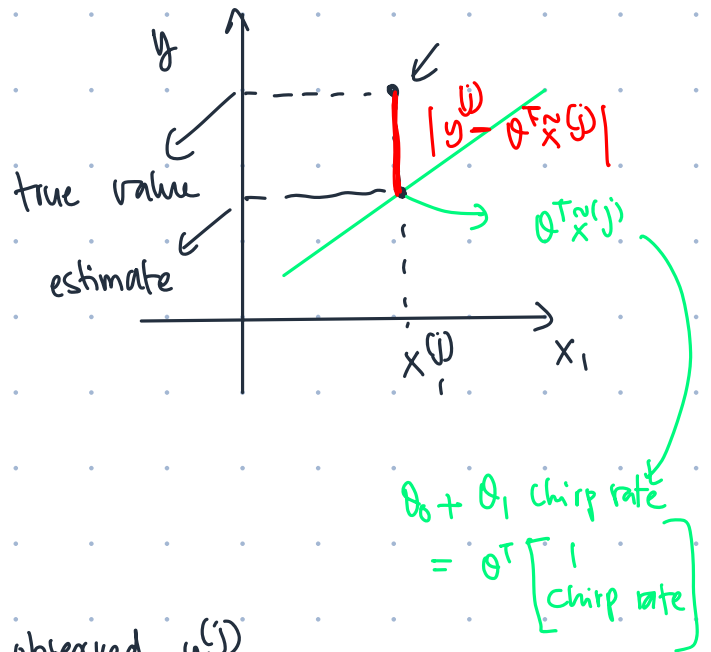
$x_0^{(j)} = 1$ ←

COST FUNCTION $J(\theta)$

GOAL: minimize
 $|y^{(j)} - \theta^T x^{(j)}|$
for all j

$$J(\theta) = \frac{1}{2n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$

average squared deviation b/w observed $y^{(j)}$,
predicted $\theta^T x^{(j)}$



How DO WE FIND "OPTIMAL" θ ? to minimize $J(\theta)$

Starting with some $\theta^{(k)}$

Nick's idea - use GD

$$\theta^{(k+1)} = \theta^{(k)} + \alpha (-\nabla J(\theta^{(k)})) \rightarrow \text{move in steepest descent direction}$$

↓
step size!

NEED : $\nabla J(\theta^{(k)})$, or more generally, $\nabla J(\theta)$

$$J(\theta) = \frac{1}{2n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}$$

$$\frac{\partial J}{\partial \theta_i} = \frac{\partial}{\partial \theta_i} \frac{1}{2n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$
$$= \frac{1}{2n} \sum_{j=1}^n \frac{\partial}{\partial \theta_i} (y^{(j)} - \theta^T x^{(j)})^2$$

$$= \sum (y^{(j)} - \theta^T x^{(j)}) \underbrace{\frac{\partial}{\partial \theta_i} (y^{(j)} - \theta^T x^{(j)})}_{-\frac{\partial \theta^T x^{(j)}}{\partial \theta_i}}$$

$$\frac{\partial}{\partial \theta_i} \theta^T x^{(j)} = \theta_0 + \theta_1 x_1^{(j)} + \dots + \theta_i x_i^{(j)} + \dots + \theta_d x_d^{(j)}$$
$$= x_i^{(j)}$$

$$\frac{\partial J}{\partial \theta_i} = -\frac{1}{n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)}) x_i^{(j)}$$

$$\theta_i^{(k+1)} = \theta_i^{(k)} + \alpha \left(\frac{1}{n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)}) x_i^{(j)} \right) \rightarrow \text{for one } \theta_i$$

$$\Rightarrow \theta^{(k+1)} = \theta^{(k)} + \frac{\alpha}{n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)}) x^{(j)}$$

GD might not converge if $\alpha > 1$ optima
to global optima

$$J(\theta) = \frac{1}{2n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2 \quad \text{is "nice"}$$

Convex, ONE minima, GD will converge!

Q. What is the cost of running one update ($k \rightarrow k+1$)
of GD, in terms of $n = \text{dataset size}$,
 $d = \text{feature dimension}$!

$$\theta^{(k+1)} = \theta^{(k)} + \frac{\alpha}{n} \sum_{j=1}^n (y^{(j)} - \theta^{(k)T} x^{(j)}) x^{(j)} \rightarrow O(d)$$

Adityan says " $O(nd)$ "

$O(nd)$

\Rightarrow seasons change, GD No update :C

Computational complexity of one update step of GD

Q: CAN WE DO BETTER?

$$\Rightarrow \theta^{(k+1)} = \theta^{(k)} + \frac{\alpha}{n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)}) x^{(j)}$$

the problem!!!

It's OK to take an "approximate" step, so long as it's faster!

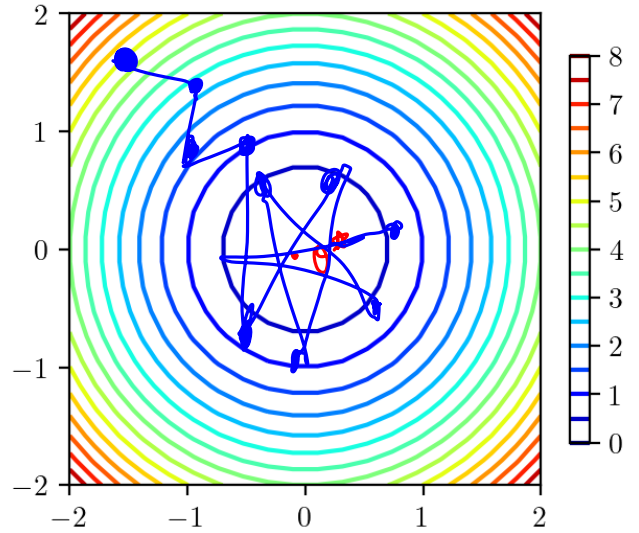
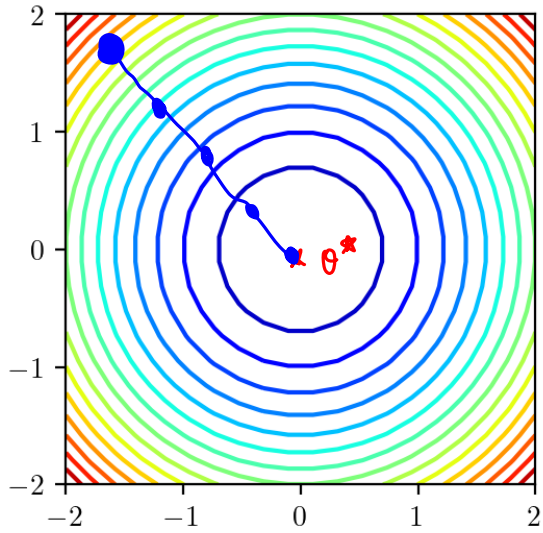
$$\nabla J = \frac{1}{n} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)}) x^{(j)}$$

1. Sample some \tilde{j} uniformly at random from training set
2. Compute $\nabla_{\tilde{j}} J$ as $(y^{(\tilde{j})} - \theta^T x^{(\tilde{j})}) x^{(\tilde{j})}$

Jay says "O(d)" for this new approach!

"Stochastic" GD

GD vs. "stochastic" GD



"Stochastic"

Another idea: closed-form solution to "least-squares" objective

"Alex" — take derivative of J w.r.t θ set 0, solve for θ ?

design matrix

$$X = \begin{bmatrix} \text{--- } x^{(1)T} \text{---} \\ \vdots \\ \text{--- } x^{(n)T} \text{---} \end{bmatrix}$$

$x^{(1)T}$ — d -dimensional

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$n \times d+1$

$$y - X\theta = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} \text{--- } x^{(1)T} \text{---} \\ \vdots \\ \text{--- } x^{(n)T} \text{---} \end{bmatrix} \theta = \begin{bmatrix} y^{(1)} - x^{(1)T}\theta \\ \vdots \\ y^{(n)} - x^{(n)T}\theta \end{bmatrix}$$

$$\|u\|_2^2 = u^T u = u_1^2 + \dots + u_n^2$$

$$\begin{aligned} \frac{1}{2n} (y - X\theta)^T (y - X\theta) &= \frac{1}{2n} \left[(y^{(1)} - x^{(1)T}\theta)^2 + \dots + (y^{(n)} - x^{(n)T}\theta)^2 \right] \\ &= \frac{1}{2n} \sum_{j=1}^n (y^{(j)} - x^{(j)T}\theta)^2 = J(\theta) \end{aligned}$$

$$J(\theta) = \frac{1}{2n} (y - X\theta)^T (y - X\theta) = (y^T - \theta^T X^T) (y - X\theta)$$

ignore for now

$$= \cancel{y^T y} - \theta^T X^T y - y^T X\theta + \theta^T X^T X\theta$$

$$\nabla_{\theta} J \stackrel{\text{set}}{=} 0 \Rightarrow \nabla_{\theta} J = \nabla_{\theta} (-\theta^T X^T y - y^T X \theta + \theta^T X^T X \theta)$$

$$a \cdot b = b \cdot a = a^T b = b^T a$$

$$\theta^T X^T y = y^T X \theta$$

$$\nabla_{\theta} (-2\theta^T X^T y + \theta^T X^T X \theta)$$

$$= -2X^T y + (X^T X + (X^T X)^T) \theta$$

$$= -2X^T y + 2X^T X \theta$$

$$\nabla_{\theta} \theta^T x = x$$

$$\nabla_{\theta} \theta^T A \theta = (A + A^T) \theta$$

$$\nabla_{\theta} J \stackrel{\text{set}}{=} 0 \rightarrow -2X^T y + 2X^T X \theta = 0$$

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

Q: Why choose GD / SGD if we have closed-form solution?

Inverting $X^T X$ is problematic — can be near singular or singular

$$\left. \begin{array}{l} (X^T X)^{-1} \rightarrow O(d^3) \\ \text{forming } X^T X \rightarrow O(nd^2) \end{array} \right\} \theta^* \rightarrow O(nd^2 + d^3)$$

vs SGD = $O(d)$

Revisiting least squares $J(\theta)$ — probabilistic view

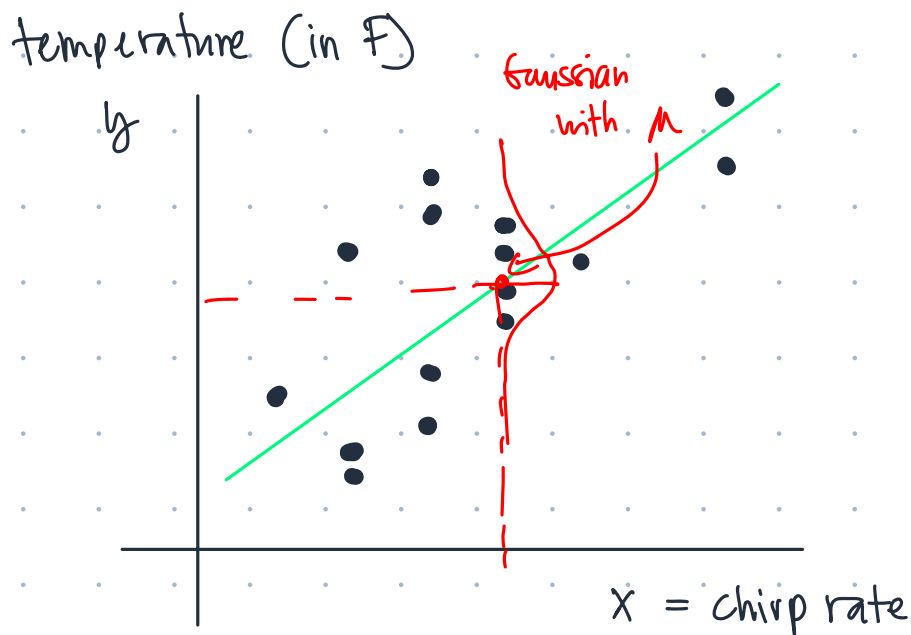
How was data generated?

$$y^{(j)} = \theta^T x^{(j)} + \epsilon^{(j)}$$

Assume: $\epsilon^{(j)} \sim \mathcal{N}(0, \sigma^2)$

$$P(\epsilon^{(j)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(j)} - 0)^2}{2\sigma^2}\right)$$

$$P(y^{(j)} | x^{(j)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(j)} - \theta^T x^{(j)})^2}{2\sigma^2}\right)$$



GOAL: Estimate " θ " to maximize likelihood of data

$$L(\theta; X, y) = P(D; \theta) = \prod_{j=1}^n P(x^{(j)}, y^{(j)}; \theta)$$

$$= \prod_{j=1}^n P(y^{(j)} | x^{(j)}; \theta) \cdot P(x^{(j)})$$

Constant

under optimization

$$\equiv \prod_{j=1}^n P(y^{(j)} | x^{(j)}; \theta)$$

$$L(\theta) = \arg \max_{\theta} \prod P(y^{(j)} | x^{(j)}; \theta)$$

$\rightarrow \mathcal{N}(\theta^T x^{(j)}; \sigma^2)$

$$l(\theta) = \log(L(\theta)) = \arg \max_{\theta} \sum_{j=1}^n \log P(y^{(j)} | x^{(j)}; \theta)$$

$$= \sum_{j=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left(-\frac{(y^{(j)} - \theta^T x^{(j)})^2}{2\sigma^2} \right) \right)$$

$$= \sum_{j=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{-(y^{(j)} - \theta^T x^{(j)})^2}{2\sigma^2}$$

constant

$$\arg \max_{\theta} \frac{1}{2\sigma^2} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$

$$\equiv \arg \min_{\theta} \frac{1}{2\sigma^2} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$

$$\equiv \arg \min_{\theta} \frac{1}{2} \sum_{j=1}^n (y^{(j)} - \theta^T x^{(j)})^2$$

Cost function,
 $J(\theta)$
 (factor of $1/n$)

Also, doesn't depend on σ^2 !