

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$

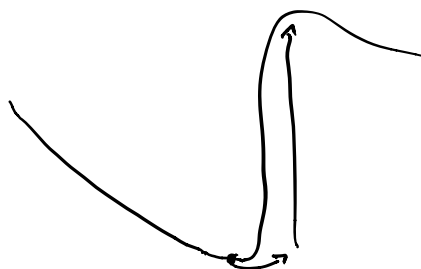
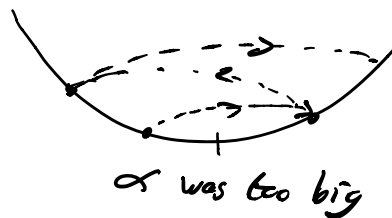
$$x^{(t+1)} = x^{(t)} - \alpha \nabla f(x^{(t)})$$

Can $f(x^{(t+1)}) > f(x^{(t)})$?

If $\alpha \gg \frac{1}{M}$, \Rightarrow divergence

$$x^{(t+1)} = (1 - \alpha M) x^{(t)}$$

$$\alpha > \frac{2}{M}$$



What morally is gradient descent?

$$\min_x f(x)$$

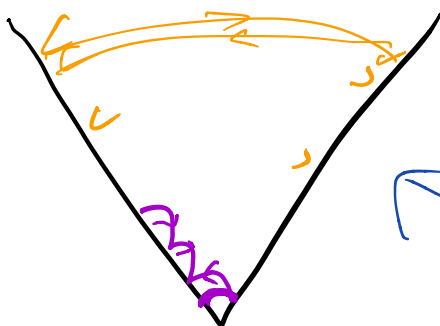
$$x^{(t+1)} = x^{(t)} - \alpha \nabla f(x^{(t)})$$

Continuous variant is gradient flow

$$\frac{dx}{dt} = - \nabla f(x(t)) \leftarrow$$

GD on $f(x) = \|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$

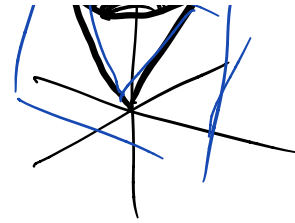
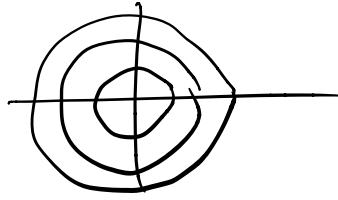
with fixed step size?



— α too big
— α not too big



$$f(x) = \sqrt{x_1^2 + x_2^2}$$



Virtues of SGD

- Can't see all data at once (online learning)
- computationally cheaper to compute approx gradient per update cost is lower
- Regularization (Stochasticity can lead us to jump out of minima)



$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i)$$

$$\frac{1}{|B|} \sum_{i \in B} \ell(f_{\theta}(x_i), y_i)$$

How big to choose B?

Largest we can afford w/ our GPU

$\{1, 2, \dots, n\}$

\uparrow
 Π

$$G(\theta) = \nabla_{\theta} \ell(f_{\theta}(x_j), y_j) \text{ for random } j$$

$$\mathbb{E}_j G(\theta) = \sum_j \frac{1}{n} \nabla_{\theta} \ell(f_{\theta}(x_j), y_j)$$

EM & Estimating Gaussian Mixtures

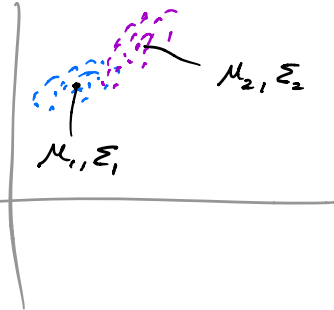
Data: $\{X_i\}_{i=1 \dots n}$

Model: Mixture of 2 Gaussians

latent variable $Z_i = \begin{cases} 1 & \text{w/ prob } \tau_1 \\ 2 & \text{w/ prob } \tau_2 = 1 - \tau_1 \end{cases}$

$$X_i | Z_i = 1 \sim \mathcal{N}(\mu_1, \Sigma_1)$$

$$X_i | Z_i = 2 \sim \mathcal{N}(\mu_2, \Sigma_2)$$



Given $\{X_i\}$, find: $\theta = \{\tau, \mu_1, \Sigma_1, \mu_2, \Sigma_2\}$

Likelihood of data under the model

$$L(\theta; X, Z) = P(X, Z | \theta) = \prod_{i=1}^n \prod_{j=1}^2 \underbrace{f(X_i; \mu_j, \Sigma_j)}_{P(X|Z, \theta)} \tau_j \underbrace{\mathbb{1}_{Z_i=j}}_{P(Z|\theta)}$$

Likelihood under normal dist

Challenge: don't know Z

① Estimate dist over Z given $\theta = \hat{\theta}$

E-step
Expectation

$$P(Z_i = j | X_i, \theta) = \frac{P(X_i | Z_i = j, \theta) P(Z_i = j | \theta)}{P(X_i | \theta)}$$

$$= \frac{f(X_i; \mu_j, \Sigma_j) \tau_j}{f(X_i; \mu_1, \Sigma_1) \tau_1 + f(X_i; \mu_2, \Sigma_2) \tau_2}$$

We can rewrite log likelihood function:

$$Q(\theta | \hat{\theta}) = \mathbb{E}_{Z | X, \hat{\theta}} \log L(\theta; X, Z)$$

$$= \sum_{i=1}^n \mathbb{E}_{Z_i | X_i, \hat{\theta}} \log L(\theta; X_i, Z_i)$$

$$= \sum_{i=1}^n \sum_j P(Z_i = j | X_i, \hat{\theta}) \log L(\theta; X_i, Z_i)$$

M-step
Maximization

② $\theta_{\text{new}} = \text{argmax } Q(\theta | \hat{\theta})$

Crowd Sourcing Problem

	P_1	P_2	P_3	P_4	...	P_m	
X_1	\hat{y}		\hat{y}		\hat{y}	\hat{y}	$\rightarrow y_1$
X_2		\hat{y}	\hat{y}		\hat{y}	\hat{y}	$\rightarrow y_2$
X_3							
X_4	\hat{y}		\hat{y}				
\vdots		\hat{y}					
X_n	\hat{y}			\hat{y}	\hat{y}		$\rightarrow y_n$

Why can EM beat MV?

Estimate quality of each voter

True class	i	P_1	P_2	P_3	P_4	P_5
(A)	1	A	A	C	C	C
(B)	2	B	B	D		
(C)	3	C	C		D	
(A)	4	A	A			B
S						

If each voter rated 1 image only?